

**Clasificador de máquinas de vectores de soporte para problemas
desbalanceados con selección automática de parámetros**

*Este trabajo de grado es presentado como un requerimiento para obtener el
título de:
Magister en Ingeniería Eléctrica*

Cristian Alfonso Jiménez Castaño

Tutor: Andrés Marino Álvarez Meza - MEng, PhD



**Universidad Tecnológica de Pereira
Facultad de Ingenierías Eléctrica, Electrónica, Física y Ciencias de la
Computación
Maestría en Ingeniería Eléctrica - Línea de automática
Grupo de Investigación Automática
Pereira-Risaralda
2018**

Índice general

1. Agradecimientos	II
2. Resumen	III
3. Abstract	IV
4. Notación	V
5. Introducción	1
5.1. Motivación	1
5.2. Planteamiento del problema	1
5.3. Estado del Arte	2
6. Objetivos	4
6.1. Objetivos General	4
6.2. Objetivos Específicos	4
I. Preliminares	5
7. Preliminares Matemáticos	6
7.1. RKHS	6
7.2. RKHS en aprendizaje de máquina	8
8. Multiplicadores de Lagrange	10
9. Preliminares Clasificación	14
9.1. Aprendizaje supervisado, clasificación desbalanceadas	14
9.2. SVM	15
9.2.1. SVM en RKHS	19
9.2.2. WSVM	20
9.3. TWSVM	21
9.3.1. TWSVM lineal	22
9.3.2. TWSVM No-lineal	25
9.4. TWSVM Delimitadas (TBSVM)	25
9.4.1. TBSVM lineal	26
9.4.2. TBSVM No-Lineal	27

9.5. WLTSVM	27
9.5.1. WLTSVM lineal	28
9.5.2. WLTSVM No-Lineal	31
II. Propuesta	32
10. Máquinas de Vectores de Soporte Gemelas mejorada (ETWSVM)	33
10.1. Fundamentos en ETWSVM	33
10.2. Aprendizaje de la función kernel en ETWSVM usando alineación centrada	36
10.3. ETWSVM para problemas Multi-clase	37
10.3.1. ETWSVM One-versus-Rest	37
10.3.2. ETWSVM One-versus-One	38
11. Montaje experimental	39
11.1. Bases de Datos	39
11.2. Entrenamiento de ETWSVM y métodos de comparación	40
11.2.1. Pre-procesamiento	40
11.2.2. Aprendizaje de la función kernel	40
11.2.3. Evaluación y sintonización de ETWSVM	41
11.2.4. Métodos de comparación	42
III. Resultados, Conclusiones y Trabajos Futuros	43
12. Resultados y Discusión	44
12.1. Resultados en datos sintéticos	44
12.2. Resultados en bases de datos reales: clasificación binaria	48
12.3. Resultados bases de datos reales: clasificación multiclase	51
13. Conclusiones	54
14. Trabajos futuros	55
15. Publicaciones	56
Bibliografía	57

Índice de figuras

7.1. Mapeo basado en función kernel.	8
8.1. Izquierda: Representación geométrica de los Multiplicadores de Lagrange	11
8.2. Representación de los multiplicadores de Lagrange	12
9.1. Ejemplo de clasificación	14
9.2. Ilustración geométrica de la distancia ortogonal de una muestra x , representada en color rojo y la superficie de decisión, vista en azul; su vector perpendicular w es ilustrado en verde.	15
9.3. La SVM busca maximizar la margen, la cual se define como la mínima distancia perpendicular del hiperplano a cualquiera de las muestras disponibles. (a) muestra un ejemplo gráfico sobre la definición del la <i>magen</i> . (b) es un ejemplo donde existe traslape entre las clases y se permite un hiperplano que clasifique muestras de una clase que no se encuentren en zona correspondiente, para ello se hace uso de las variables auxiliares ξ_n . En ambas figuras, las muestras encerradas en un círculo verde son consideradas vectores de soporte y son las que determinan el hiperplano.	16
9.4. La TWSVM genera dos hiperplanos no-paralelos, donde cada uno esta dedicado a una clase. Cada hiperplano debe estar lo más cerca a su clase correspondiente y lo más alejado posible de las muestras de la otra clase. Las muestras e hiperplanos de color rojo corresponden a la clase $+1$ y los de color azul a la clase -1 . Las muestras que se encuentran cerradas en las circunferencias verdes, son las muestras que definen los márgenes de los hiperplanos.	23
9.5. a) Densidad intraclase según el submuestreo basado en ν -nn. b) Densidad interclase según el submuestreo basado en ν' -nn. Para la generación de ambas gráficas fijamos $\nu=\nu'=6$	29
10.1. Representación geométrica de las estrategias multiclase acopladas con el ETWSVM. a) corresponde al esquema OvR acoplado con el ETWSVM, donde se generan $R=3$ hiperplano, uno por cada clase. b) representa el esquema OvO acoplado con el ETWSVM, donde para $R=3$, se general tres sub-clasificadores (por cada uno dos hiperplanos).	38
11.1. Disposición del ETWSVM*. Un submuestreo basado en ν -NN es aplicado antes de la etapa del aprendizaje de la función kernel. Entonces, la clasificación binaria o multiclase es realizado para resolver el problema de optimización de ETWSVM.	40

12.1. Resultados de la clasificación del ETWSVM* sobre la base de datos media-luna. (a) muestra los datos sintéticos de media-luna, las muestras azules provienen de la clase minoritaria mientras las muestras rojas de la clase mayoritaria; las instancias encerradas en círculos verdes son seleccionadas por el submuestreo basado en ν -NN. (b) muestra la proyección basada en CKA fijando $P'=2$. (c) y (d) muestra la distancia entre los datos de entrada y los hiperplanos minoritario y mayoritaria, respectivamente; donde las muestras encerradas en negro y en cian representa los vectores de soporte de los hiperplanos minoritario y mayoritario, respectivamente. (e) muestra la frontera de decisión y la mínima distancia entre (c) y (d) (ver función de decisión del ETWSVM en ecuación (10.24)).	44
12.2. Inspección visual de los resultados sobre bases de datos sintéticas generada a partir de dos distribuciones Gaussianas multivariadas (clasificación basada en ETWSVM*). Diferentes escenarios de desbalance y traslape entre las clases son probadas. Las muestras azules provienen de la clase minoritaria mientras que las rojas pertenecen a la mayoritaria. Las muestras encerradas en negro y cian representan los vectores de soporte de las clase minoritaria y mayoritaria, respectivamente, y la cruz la media de las distribuciones. La frontera de decisión es representada con una línea negra y la mínima distancia a los hiperplanos es coloreada en el fondo.	45
12.3. Análisis de los parámetros de regularización (resultados de clasificación del ETWSVM*). La evaluación cualitativa basada en el GM se representa variando la tasa de desbalance y la distancia entre las medias de cada clase. Dentro de cada gráfica un par de valores de $c_{1,\ell}$ y $c_{2,\ell}$ es fijado para solucionar los QPPs en la ecuación (10.20). Esta medida de rendimiento se calcula variando de abajo hacia arriba la relación de desbalance y de izquierda a derecha la distancia entre las clases.	46
12.4. Análisis de los parámetros de regularización (resultados de clasificación del ETWSVM*). La evaluación cualitativa basada en el FM se representa variando la tasa de desbalance y la distancia entre las medias de cada clase. Dentro de cada gráfica un par de valores de $c_{1,\ell}$ y $c_{2,\ell}$ es fijado para solucionar los QPPs en la ecuación (10.20). Esta medida de rendimiento se calcula variando de abajo hacia arriba la relación de desbalance y de izquierda a derecha la distancia entre las clases.	47
12.5. Frontera de decisión de la base de datos Vehicle después de la embebimiento a un espacio 2D a través del t-SNE. Las muestras azules provienen de la clase minoritaria mientras que las rojas pertenecen a la clase mayoritaria. Las fronteras de decisión son mostradas en negro, y el puntaje de la clasificación (o la mínima distancia entre a los hiperplanos) son coloreados en el fondo.	49
12.6. Prueba de diferencia estadística entre los clasificadores basada en Tukey-Kramer. Se muestra el diagrama de caja del ranking de media del algoritmo, y el color indica que algoritmos son estadísticamente iguales.	51
12.7. Prueba de significancia estadística entre los clasificadores estudiados sobre el repositorio Keel, usando la prueba de Tukey-kramer. Se muestra el diagrama de caja del ranking de media del algoritmo, y el color indica que algoritmos son estadísticamente iguales.	53

Índice de tablas

11.1. Repositorios UCI y Keel descritos para clasificación desbalanceada. $IR = N_+/N_-$: tasa de desbalance, P : # de características, N : # de muestras, y R : # de clases.	39
12.1. Resultados sobre las bases de datos del repositorio UCI (clasificación binaria). Las medidas Acc, GM y FM son consideradas. La media \pm la desviación estándar son mostradas como resultado de la validación cruzada anidada de 10-fold.	50
12.2. Tiempo de entrenamiento para las bases de datos del repositorio UCI (clasificación binaria). La media \pm desviación estándar son representan el tiempo de entrenamiento en segundos por cada fold.	50
12.3. Resultados sobre las bases de datos del repositorio Keel (clasificación multiclase). Las medidas \overline{Acc} y el FM_μ son considerados. La media \pm desviación estándar son mostradas para la validación cruzada anidada de 10-fold.	52
12.4. Tiempo de entrenamiento para el repositorio Keel (clasificación multiclase). La media \pm la desviación estándar mostradas corresponden al tiempo de entrenamiento y esta dado en segundos por fold.	52

1. Agradecimientos

Quiero agradecer a mi familia por su continuo e incondicional apoyo durante toda mi vida, en especial a mis padres, que nunca han dejado de enseñarme. También en especial a mi tío, mi segundo padre, José Nelson por ese apoyo incondicional a mí y a toda la familia, por demostrarnos que con trabajo arduo, dedicación y motivación todo es posible. A mis Abuelas y Abuelo por ser esas personas incondicionales en mi vida, que con amor y comprensión nos han ayudado a pararnos y seguir adelante. También, quiero agradecerle a mi tutor Andrés Marino Álvarez M, por sus contribuciones a mi trabajo, a mi formación como persona y profesional, y en especial por esa calidad humana que posee. Estoy más que agradecido con la Maestría en Ingeniería Eléctrica del programa de Ingeniería Eléctrica en la Universidad Tecnológica de Pereira y con el grupo de investigación en Automática.

Esta investigación es desarrollada bajo el proyecto “Desarrollo de un sistema de soporte clínico basado en el procesamiento estocástico para mejorar la resolución espacial de la resonancia magnética estructural y de difusión con aplicación al procedimiento de la ablación de tumores”, code: 111074455860, financiado por COLCIENCIAS. Además este trabajo fue parcialmente financiado por el proyecto E6-18-09: “Clasificador de máquinas de vectores de soporte para problemas desbalanceados con selección automática de parámetros”, de la Vicerrectoría de Investigación, innovación y extensión de la Universidad Tecnológica de Pereira.

2. Resumen

La mayoría de los métodos de clasificación asumen que el número de muestras en las clases estudiadas son las mismas (balanceadas). Sin embargo, realizar esta asunción puede llevar a desempeños sesgados, ya que, la mayoría de aplicaciones y bases de datos reales no son balanceadas, llevando a que estos métodos ignoren la clase minoritaria (la clase con el menor número de muestras). Este trabajo propone un clasificador novedoso, llamado *enhanced twin support vector machine*–(ETWSVM), que representa las muestras de entrada en un espacio de características de alta dimensionalidad, posiblemente infinita, durante la construcción de una frontera de decisión bajo la filosofía del *twin support vector machine*–(TWSVM). También, usamos un método basado en *centered kernel alignment*–(CKA) para aprender la función kernel con el fin de contrarrestar los problemas inherentes del desbalance y mejorar la separabilidad de los datos. Además, adoptamos las estrategias *One-versus-Rest* y *One-versus-One* para extender la formulación del ETWSVM a tareas de clasificación multiclase. De los resultados obtenidos sobre bases de datos sintéticas y reales, nuestra propuesta supera métodos del estado del arte con respecto al desempeño (precisión, media geométrica, F-measure), y tiempo de entrenamiento. En efecto, después analizamos la sensibilidad de los parámetros libres para diferentes tasas de desbalance y traslape entre las clases, y sugerimos una variante del ETWSVMN automático que registra una indicada relación entre desempeño de clasificación y tiempo de entrenamiento.

3. Abstract

Most of the classification approaches are based on the assumption that the sample distribution among classes is balanced. Nonetheless, such an assumption leads to biased performance over the majority class (the class with the highest number of samples). This work proposes a novel approach, called *enhanced twin support vector machine*–(ETWSVM), that represents the input samples in a high-dimensional feature space of possible infinity dimension during the decision boundary building under a *twin support vector machine*–TWSVM philosophy. Also, we use a *centered kernel alignment*–(CKA)-based approach to learning the kernel function counteracting inherent imbalanced issues and improving the data separability. Moreover, we also adopt One-vs-One and One-vs-Rest frameworks to extend the ETWSVM formulation for multi-class tasks. Obtained results on synthetic and real-world datasets show that our approach outperforms state-of-the-art methods concerning classification performance (accuracy, geometric mean, and F-measure), and training time. Indeed, after analyzing the sensitive of the free parameters for the imbalanced ratio and the overlapping among clusters, we suggest an automatic ETWSVM variant that achieves a suitable trade-off between classification performance and training time.

4. Notación

Símbolo	Descripción
\mathbb{R}^P	espacio Euclidiano de dimensión P
a	los escalares se representaran por letras minúsculas
P	los escalares que indican dimensiones y tamaño de vectores y/o matrices se expresarán en mayúscula y sin negrilla. En especial, este indica la dimensión del espacio de entrada de los datos
P'	dimensión del espacio embebido por selección o extracción de características
Q	dimensión del espacio de características, RKHS, \mathcal{H}
N_+	número muestras de entrada pertenecientes a la clase minoritaria
N_-	número muestras de entrada pertenecientes a la clase mayoritaria
N	número total de muestras de entradas $N=N_+ + N_-$
N_S	Número de vectores de soporte total del modelo entrenado
ℓ	índice que hace referencia a la clase mayoritaria (-1) o minoritaria ($+1$), por facilidad este toma valores del conjunto $\ell \in \{+, -\}$. Su reciproco es $\ell' = -\ell$
\mathbf{x}	los vectores columnas serán representados en letras minúsculas y en negrilla. Se asumirá que estos pertenecen a \mathbb{R}^P siempre y cuando no se especifique lo contrario
\mathbf{A}	las matrices serán indicadas con letras mayúsculas y en negrilla
$a_{nn'}$	elemento de la fila n y columna n' de la matriz \mathbf{A}
$\boldsymbol{\alpha}_{\ell'}$	vector de los multiplicadores de Lagrange de las muestras de la clase $\ell'1$, donde $\ell' \in \{+, -\}$
\mathbf{I}	matriz identidad, una matriz cuadrada de dimensiones apropiadas, tiene 1's en su diagonal y 0's fuera de esta
$\ \cdot\ _2$	norma Euclidiana o también conocida como norma L_2
$\ \cdot\ _{\mathcal{H}}$	norma Euclidiana o también conocida como norma L_2 en el espacio RKHS \mathcal{H}
$ \cdot $	valor absoluto de \cdot
$\langle \mathbf{x}_i, \mathbf{x}_j \rangle_{\mathcal{H}}$	producto interior entre los vectores columna \mathbf{x}_i y \mathbf{x}_j in el espacio \mathcal{H} , donde $\mathbf{x}_i, \mathbf{x}_j \in \mathcal{H}$
$\mathbf{1}_{\ell}$	vector de longitud N_{ℓ} con sólo valores de 1's
$\mathbf{0}_M$	vector de longitud M con sólo valores de 0
$f(\cdot)$	función escalar evaluada en \cdot
$\text{sign}(\cdot)$	función <i>signo</i> evaluada en \cdot
$\text{med}(\cdot)$	operador mediana de \cdot
$\varphi(\cdot)$	función de mapeo desde el espacio de entrada \mathcal{X} al espacio de características, o RKHS, \mathcal{H}

5. Introducción

5.1. Motivación

En la formulación de métodos de aprendizaje para clasificación una de las primeras asunciones es que el número de muestras entre las clases estudiadas son iguales, es decir que se va aplicar sobre bases de datos balanceadas. El problema de desbalance de datos en clasificación ha captado, en los últimos años, un gran interés por parte de la comunidad científica. Debido, a que varias aplicaciones del mundo real tienen muchas muestras de un fenómeno en comparación al número de muestras de otro, algunos ejemplos de esto se pueden encontrar en: la detección de intrusión [1], el diagnóstico médico [2], el diagnóstico de cáncer [3], el análisis de tráfico [4], y valoración de crédito [5].

El principal problema del desbalance de datos es la gran posibilidad de que la clase minoritaria, la clase con el menor número de muestras, sea sesgada por el clasificador. Debido a que esta clase puede ser tratada como ruido, o que el ruido puedan ser tratados como de esta clase, por el clasificador [6].

5.2. Planteamiento del problema

En la actualidad este problema es atacado principalmente desde tres enfoques. El primero es a nivel de los datos en el cual se realiza un pre-procesamiento de los mismos, el más común son los métodos de remuestreo en los cuales se busca balancear el número de muestras entre las clases por medio de la eliminación y/p generación de muestras de las clases [7]. El segundo enfoque es el desarrollo de algoritmos costo-sensible, donde se pondera la función de costo del algoritmo de aprendizaje de tal forma que la mala clasificación de la clase minoritaria se penalice con mayor fuerza [8]. El tercer enfoque es el ensamblaje de clasificadores, donde se fragmenta la clase mayoritaria en subconjuntos de la misma, sin que se sobrepongan, y son entrenados tantos clasificadores como subconjuntos obtenidos, la predicción se realiza por votación mayoritaria [9].

Entendiendo un poco lo anterior. El remuestreo se puede dividir en tres categorías: submuestreo de la clase mayoritaria, sobremuestreo de la clase minoritaria y el híbrido entre ellos. El submuestreo, consiste en extraer muestras de la clase mayoritaria para así balancear el número de muestras entre ambas clases, el principal problema del submuestreo es la pérdida potencial de información debido a la eliminación de las muestras procedentes de la clase mayoritaria. Habría que decir también, que el sobremuestreo consiste en expandir el conjunto de muestras de la clase minoritaria por medio de la creación de muestras sintéticas, a partir de la originales; pero este tiende a generar sobre-entrenamiento por la información redundante procedente de estas muestras, además, el tiempo de cómputo crece drásticamente debido a la creación de dichas muestras. Cabe destacar, que los métodos de remuestreo son indiferentes al clasificador a usar [7]. El segundo enfoque, ensamblaje de clasificadores, consiste en fragmentar el conjunto de muestras mayoritarias en subconjuntos de las mismas, de tal forma que no exista intersección entre ellas, del mismo tamaño de la clase minoritaria. Luego, por cada subconjunto de la clase mayoritaria y el conjunto de la clase minorita-

ria se entrena un clasificador, al final se entrenan tantos clasificadores como subconjuntos se extraigan. Por lo general las metodologías basadas en esta filosofía presentan un buen rendimiento pero su costo computacional es alto, por el número de clasificadores a entrenar [10]. Por último, el enfoque por algoritmos de costo-sensible pueden llevar fácilmente a sobre-entrenamiento [8].

En años recientes, han derivado métodos basados en máquinas de vectores de soporte (*support vector machine*-(SVM)) con el fin de replicar su capacidad de generalización con un costo computacional menor [11]. Uno de estos métodos son las máquinas de vectores de soporte gemelas (*twin support vector machine*-(TWSVM)) las cual soluciona dos problemas de programación cuadráticas (QPPs) más pequeñas que el QPP que soluciona la SVM. Lo anterior, lleva a que las TWSVM tengan un costo computacional menor que las SVM con una buena capacidad de generalización [12]. No obstante, poseen el mismo problema ante bases de datos desbalanceadas que la SVM estándar. Algunas extensiones del TWSVM incluyen una representación de la estructura de los datos para favorecer su capacidad de generalización de la clasificación [13]. Adicionalmente, han sido propuestos algoritmos basados en TWSVM para contrarrestar los problemas de desbalance de datos, combinándolo con técnicas de re-muestreo, además de problemas de optimización ponderados con el fin de evitar sesgo en la discriminación [14–16]. Aún así, los enfoques basados en TWSVM no incluyen una representación apropiada a un espacio de Hilbert con kernel reproductivo (RKHS) por medio una función de mapeo, la cual mapea las muestras del espacio original a un espacio de alta dimensionalidad que mejore la representación de los datos. Por lo tanto, descuidan la formulación intrínseca y las virtudes del problema dual con respecto a los métodos kernels.

De lo anterior se puede concluir que, los principales problemas del desbalance de datos son: tiempo de cómputo, sobre-entrenamiento y pérdida de información. También se puede concluir que el problema de desbalance no se ha podido solucionar por completo y nos deja con la pregunta de investigación: ¿Es posible el desarrollo e implementación de una metodología basada en TWSVM con una nueva extensión no-lineal con mapeo a un espacio de alta dimensionalidad, posiblemente infinito, el cual contrarreste los efectos del desbalance en problemas de clasificación?, además ¿es posible la selección automática de los parámetros libres en aras de reducir el tiempo de entrenamiento y favorecer la identificación de patrones relevantes?.

5.3. Estado del Arte

En el estado del arte las técnicas de remuestreo se dividen en: submuestreo de la clase mayoritaria, sobre-muestreo de la clase minoritaria e híbrido. La técnica de submuestreo más eficiente y simple es el submuestreo aleatorio (RUS, por sus siglas en inglés), el cual consiste en la eliminación aleatoria de muestras de la clase mayoritaria [17]. Dos métodos ampliamente utilizados como técnicas de sobremuestreo son duplicación aleatoria de muestras minoritarias y el *synthetic minority over-sampling technique*-(SMOTE) [18]. [19] realizó un estudio de los tipos de clases minoritarias y su influencia en el aprendizaje del clasificador en datos desbalanceados. Algunos autores han tratado de definir de forma automática la proporción de muestreo para diferentes tasas de desbalance y configuración del problema [20–25]. Por otra parte, [7, 26] han realizado un estudio sobre el desempeño de diferentes técnicas de remuestreo, de los cuales se puede concluir: que cuando se tiene cientos de muestras de la clase minoritaria, el submuestreo supera al sobremuestreo en términos de tiempo de computó; cuando se tienen pocas muestras de la clase minoritaria, la mejor opción es aplicar alguna técnica de sobremuestreo; si el tamaño de las muestras de entrenamiento es muy grande, lo mejor es combinar ambos tipos de técnicas de remuestreo; para finalizar, los métodos de sobre-muestreo son un poco más eficaces para reconocer los datos atípicos.

Por otra parte, la estrategia de ensamblamiento de clasificadores es muy popular pero conlleva un gran costo

computacional para entrenar varios clasificadores. Luego, esta estrategia se divide principalmente en dos enfoques: Bagging y Boosting. El Bagging consiste en que los clasificadores base pueden ser entrenados en paralelos. Por otro lado, el Boosting consiste en un algoritmo de entrenamiento de los clasificadores bases, el cual es evolutivo y calcula los pesos que van a tener cada clasificador base en la predicción. Este último enfoque es el más común y eficiente, el primer método de esta naturaleza fue AdaBoost propuesto en 1996 [27]; en el estado del arte se encuentra una gran variedad de métodos basado en este y con diferentes tipos de clasificadores base, algunos ejemplos son: [9, 10, 28, 29].

Por último, los algoritmos de costo-sensibles comparados con los métodos de remuestreo, los algoritmos de aprendizaje son más eficientes computacionalmente hablando, por lo que puede ser más adecuado para grandes flujos de datos. Sin embargo, es más popular los métodos de remuestreo para escenarios de desbalance por su simplicidad y su independencia del clasificador. Dentro del estado del arte, se encuentran métodos donde se cambia el proceso de aprendizaje o la función de costo para construir un clasificador robusto al desbalance [30–35]. Entre los anteriores se encuentran el *weighted lagrangian twin support vector machine*–(WLTSVM) [30], el cual es una variación de la máquina de vectores de soporte gemelas (TWSVM) [36], que tuvo resultados muy buenos en cuanto al tiempo de entrenamiento.

En este trabajo, proponemos un clasificador basado en máquinas de vectores de soporte gemelas que incluye una función kernel para mejorar la formulación del mismo para clasificación de datos desbalanceados. Nuestra propuesta es denominada *enhanced TWSVM*–(ETWSVM), la cual permite la representación de los datos de entrada en un espacio de alta dimensionalidad, posiblemente infinita dimensión, durante la optimización del clasificador. También, usamos un algoritmo basado en *centered kernel alignment*–(CKA) [37, 38] con el objetivo de aprender la función kernel de mapeo del ETWSVM para contrarrestar los problemas inherentes al desbalance y mejorar la separabilidad de los datos. Además, extendemos nuestro clasificador biclase a multiclase utilizando las estrategias One-versus-Rest (OvR) y One-versus-One (OvO) [39]. Como puntos de referencia para la clasificación binaria, probamos la SVM estándar [40], la SVM con variables slack regularizadas o WSVM [41], un SMOTE-SVM (SVM_{SMOTE}) [42], el mejorado TWSVM–(TBSVM) [12], y las máquinas de vectores de soporte gemelas con Lagrangiano ponderado (WLTSVM) [30]. Para la clasificación multi-clase, probamos SVM y el TBSVM con los esquemas OvR y OvO. Los resultados obtenidos sobre las bases de datos referentes del estado del arte muestran que nuestra propuesta supera los métodos de referencia en cuanto a las medidas de rendimiento de precisión, media geométrica y F-measure (F_1 -score). También, estudiamos la influencia de los parámetros libre del ETWSVM considerando el desempeño del clasificador para diferentes tasas de desbalance y escenarios de traslape, sugerimos una variante de ETWSVM que obtiene un notable equilibrio entre desempeño y tiempo de entrenamiento, sin la necesidad de conocimiento a priori sofisticado con respecto a la sintonización del algoritmo.

El resto de este documento esta organizado como sigue: La parte I trata los temas preliminares, algoritmos de aprendizaje de gran importancia para nuestra propuesta. La parte II contiene lo relacionado a la teoría de nuestra propuesta y sus variaciones. Por último, la parte III contiene los resultados obtenidos con su respectiva discusión, además de las conclusiones y trabajos futuros.

6. Objetivos

6.1. Objetivos General

Desarrollar una metodología de clasificación basada en máquinas de vectores de soporte gemelas que considere el desbalance entre clases y que contemple la selección automática de parámetros libres en aras de favorecer la identificación de patrones relevantes con base en criterios de sensibilidad y especificidad.

6.2. Objetivos Específicos

- Proponer un método de clasificación basado en máquinas de vectores de soporte gemelas (TWSVM) que contemple mapeos no lineales con kernel reproductivo dentro de un esquema de optimización cuadrático, con el fin de favorecer la generalización de las representaciones utilizando distintos kernels según la aplicación de interés.
- Desarrollar una estrategia de sintonización de los parámetros libres del método de clasificación propuesto, que incluya la regularización automática del problema de aprendizaje y la escogencia de la función kernel a partir de los datos disponibles.
- Desarrollar una metodología de remuestreo para codificar estructuras de datos relevantes y favorecer el rendimiento en tareas de clasificación con desbalance, a partir del método de clasificación propuesto basado en TWSVM con selección automática de parámetros.

Parte I.

Preliminares

7. Preliminares Matemáticos

En este capítulo, proporcionamos una breve introducción a los conceptos de la teoría de funciones kernels para representación de muestras de entrada en el contexto de aprendizaje de máquina. Primero, se plantean las definiciones formales y necesarias, además de las condiciones suficientes para que una función sea representada por un *kernel reproductivo*. El contenido de este capítulo está basado en artículos y libros de Parzen [43], Aronszajn [44], el Scholkopf y Smola [11].

7.1. Espacios de Hilbert con kernel reproductivo

Sea \mathcal{X} un conjunto y \mathcal{F} un espacio vectorial de funciones en \mathcal{X} para el campo \mathbb{F} ; en particular, $\mathbb{F}=\mathbb{R}$. Entonces, existe un espacio de Hilbert con kernel reproductivo (*reproductive kernel Hilbert space*–(RKHS)) \mathcal{H} en \mathcal{X} sobre \mathbb{R} , si:

- \mathcal{H} es un sub-espacio vectorial de \mathcal{F} .
- \mathcal{H} está dotado de un producto punto, $\langle \cdot, \cdot \rangle_{\mathcal{H}}$, y es completo dentro de la métrica inducida por este.
- Para cada $\mathbf{x} \in \mathcal{X}$ y $f \in \mathcal{H}$, la evolución lineal funcional $F_{\mathbf{x}} : \mathcal{H} \rightarrow \mathbb{R}$, se define como $F_{\mathbf{x}}(f) = f(\mathbf{x})$, es acotada

Desde el teorema de Riesz [45], es de conocimiento que para cualquier función acotada H en un espacio de Hilbert \mathcal{H} , existe un único vector $\mathbf{h} \in \mathcal{H}$ tal que: $H(\mathbf{f}) = \langle \mathbf{h}, \mathbf{f} \rangle_{\mathcal{H}} \forall \mathbf{f} \in \mathcal{H}$. En consecuencia, para cada evaluación funcional $F_{\mathbf{x}}$ existe un correspondiente vector $\kappa_{\mathbf{x}} \in \mathcal{H}$. La función bivariable se define como

$$\kappa(\mathbf{x}, \mathbf{x}') = \kappa_{\mathbf{x}}(\mathbf{x}'), \quad (7.1)$$

esta función es llamada *kernel reproductivo* para \mathcal{H} , con $\mathbf{x}' \in \mathcal{X}$. Se puede verificar que:

$$\kappa(\mathbf{x}, \mathbf{x}') = \langle \kappa_{\mathbf{x}}, \kappa_{\mathbf{x}'} \rangle_{\mathcal{H}}, \quad (7.2)$$

y $\|F_{\mathbf{x}}\|_{\mathcal{H}}^2 = \|\kappa_{\mathbf{x}}\|_{\mathcal{H}}^2 = \langle \kappa_{\mathbf{x}}, \kappa_{\mathbf{x}} \rangle_{\mathcal{H}} = \kappa(\mathbf{x}, \mathbf{x})$, donde $\|\cdot\|$ es el operador norma.

Sea \mathcal{H} un RKHS en el conjunto \mathcal{X} con kernel $\kappa(\cdot, \cdot)$. El tramo lineal de $\{\kappa(\mathbf{x}, \cdot) : \mathbf{x} \in \mathcal{X}\}$ es denso en \mathcal{H} . Este resultado del hecho de que cualquier función f ortogonal para el tramo de $\{\kappa(\mathbf{x}, \cdot) : \mathbf{x} \in \mathcal{X}\}$ debe satisfacer $\langle f, \kappa_{\mathbf{x}} \rangle_{\mathcal{H}} = 0$, y así $f(\mathbf{x}) = 0$.

Lema 7.1. Sea $\{f_n\} \subset \mathcal{H}$, con $n \in \mathbb{N}$ como un índice contable. Si $\lim_{n \rightarrow +\infty} \|f_n - f\|_{\mathcal{H}} = 0$, entonces $f(\mathbf{x}) = \lim_{n \rightarrow +\infty} f_n(\mathbf{x})$ para todo $\mathbf{x} \in \mathcal{X}$.

Demostración. Esto es una simple consecuencia de la propiedad reproductiva y la desigualdad de Cauchy-Schwarz:

$$|f_n(\mathbf{x}) - f(\mathbf{x})| = |\langle f_n - f, \kappa_{\mathbf{x}} \rangle_{\mathcal{H}}| \leq \|f_n - f\|_{\mathcal{H}} \|\kappa_{\mathbf{x}}\|_{\mathcal{H}} \rightarrow 0$$

□

Proposición 7.2. Sean \mathcal{H}_1 y \mathcal{H}_2 RKHS en \mathcal{X} con funciones kernels κ_1 y κ_2 , respectivamente. Si $\kappa_1(\mathbf{x}, \mathbf{x}') = \kappa_2(\mathbf{x}, \mathbf{x}')$ para todo $\mathbf{x}, \mathbf{x}' \in \mathcal{X}$, entonces $\mathcal{H}_1 = \mathcal{H}_2$ y $\|f\|_{\mathcal{H}_1} = \|f\|_{\mathcal{H}_2} \forall f$.

Demostración. Podemos tomar $\kappa(\mathbf{x}, \mathbf{x}') = \kappa_1(\mathbf{x}, \mathbf{x}') = \kappa_2(\mathbf{x}, \mathbf{x}')$ y así el $M_l = \text{span}\{\kappa_{\mathbf{x}} \in M_l : \mathbf{x} \in \mathcal{X}\}$ es denso en \mathcal{H}_l , y para cualquier $f(\mathbf{x}) = \sum_n \alpha_n \kappa_{\mathbf{x}_n}$ no existe consideración sobre si f pertenece a M_1 o a M_2 . Note que $\|f\|_{\mathcal{H}_1}^2 = \sum_{n,n'} \alpha_n \alpha_{n'} \kappa(\mathbf{x}_n, \mathbf{x}_{n'}) = \|f\|_{\mathcal{H}_2}^2$ y así $\|f\|_{\mathcal{H}_1} = \|f\|_{\mathcal{H}_2}$ para todo $f \in M_1 = M_2$. Si $f \in \mathcal{H}_1$, entonces existe una secuencia de funciones $\{f_n\} \subset M_1$ que converge a f en norma. Como, $\{f_n\}$ es Cauchy en M_1 es también Cauchy en M_2 , así por lo completo de \mathcal{H}_2 existe $g \in \mathcal{H}_2$ tal que $f_n \rightarrow g$. Entonces, por el lema 7.1 tenemos que $f(\mathbf{x}) = \lim_{n \rightarrow +\infty} f_n(\mathbf{x}) = g(\mathbf{x})$ para todo $\mathbf{x} \in \mathcal{X}$, cada $f \in \mathcal{H}_1$ esta también en \mathcal{H}_2 y viceversa, y por ende $\mathcal{H}_1 = \mathcal{H}_2$. Finalmente, podemos extender que $\|f\|_{\mathcal{H}_1} = \|f\|_{\mathcal{H}_2}$. □

Por lo tanto, dos RKHS diferentes no tienen el mismo kernel reproductivo. El siguiente teorema muestra un camino alternativo para expresar el kernel reproductivo de un RKHS \mathcal{H} .

Teorema 7.3. Sea \mathcal{H} un RKHS con kernel reproductivo κ . Si $\{e_\lambda(\mathbf{x}) : \lambda \in \Lambda\}$ es una base ortonormal de \mathcal{H} , entonces:

$$\kappa(\mathbf{x}, \mathbf{x}') = \sum_{\lambda \in \Lambda} e_\lambda(\mathbf{x}) e_\lambda(\mathbf{x}'), \quad (7.3)$$

donde las series convergen de forma puntual.

Demostración. Para un conjunto fijo $\{\mathbf{x}_n\} \subseteq \mathcal{X}$, tenemos:

$$\sum_{n,n'=1}^N \alpha_n \alpha_{n'} \kappa(\mathbf{x}_n, \mathbf{x}_{n'}) = \left\langle \sum_n \alpha_n \kappa_{\mathbf{x}_n}, \sum_{n'} \alpha_{n'} \kappa_{\mathbf{x}_{n'}} \right\rangle_{\mathcal{H}} = \left\| \sum_n \alpha_n \kappa_{\mathbf{x}_n} \right\|_{\mathcal{H}}^2 \geq 0$$

□

Adicional a esto, el teorema de Moore es introducido, el cual es contrario al anterior resultado y nos proporciona la característica de una función definida positiva como condición suficiente para que la función sea un kernel reproductivo de algún RKHS.

Teorema 7.4. Sea \mathcal{X} un conjunto y $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ una función definida positiva. Entonces, existe un RKHS \mathcal{H} de funciones en \mathcal{X} , tal que, κ es un kernel reproductivo de \mathcal{H} .

Demostración. Considere las funciones $\kappa_{\mathbf{x}}(\mathbf{x}') = \kappa(\mathbf{x}, \mathbf{x}')$ y el espacio W extendido por el conjunto $\{\kappa_{\mathbf{x}} : \mathbf{x} \in \mathcal{X}\}$. El siguiente mapeo bilineal $B : W \times W \rightarrow \mathbb{R}$:

$$B \left(\sum_n \alpha_n \kappa_{\mathbf{x}_n}, \sum_{n'} \beta_{n'} \kappa_{\mathbf{x}_{n'}} \right) = \sum_{n,n'} \alpha_n \beta_{n'} \kappa(\mathbf{x}_n, \mathbf{x}_{n'}),$$

donde $\alpha_n, \beta_{n'} \in \mathbb{R}$, esta bien definido en W . Para respaldar lo dicho anteriormente, notese que si $f(\mathbf{x}) = \sum_n \alpha_n \kappa_{\mathbf{x}_n}(\mathbf{x})$ es cero para todo $\mathbf{x} \in \mathcal{X}$, entonces por definición $B(f, \kappa_{\mathbf{x}}) = 0 \forall \mathbf{x}$. A la inversa, si $B(f, w) = 0$ para todo $w \in W$, entonces tomando $w = \kappa_{\mathbf{x}}$ se puede ver que $f(\mathbf{x}) = 0$. Entonces, B esta bien definida.

Como κ es definida positiva $B(f, f) \geq 0$ y vemos que $B(f, f)=0$ si y sólo si $B(w, f)=0$ para todo $w \in W$, por lo tanto $f(x)=0$ para todo $x \in \mathcal{X}$. Ahora, hemos visto que W es un sub-espacio pre-Hilbert con producto interno B . Dejemos que \mathcal{H} denote la terminación de W , necesitamos ver que cada elemento de \mathcal{H} es función sobre \mathcal{X} . Sea $h \in \mathcal{H}$ el punto limite de una secuencia de Cauchy $\{f_n\} \subseteq W$. Por la desigualdad de Cauchy-Shwarz:

$$|f_n(x) - f_{n'}(x)| = |B(f_n - f_{n'}, \kappa_x)| \leq \|f_n - f_{n'}\| \kappa(x, x).$$

Por lo tanto, el limite puntual $h(x) = \lim_{n \rightarrow +\infty} f_n(x)$ esta bien definida. Concluyendo, sea $\langle \cdot, \cdot \rangle_{\mathcal{H}}$ el producto interior sobre \mathcal{H} . Entonces, tenemos $\langle h, \kappa_x \rangle_{\mathcal{H}} = \lim_{n \rightarrow +\infty} \langle f_n, \kappa_x \rangle_{\mathcal{H}} = B(f_n, \kappa_x) = h(x)$. Así \mathcal{H} es un RKHS con kernel reproductivo κ . \square

Combinando la proposición 7.2 con el teorema de Moore (teorema 7.4) vemos la correspondencia entre RKHS sobre el conjunto \mathcal{X} y funciones definidas positivas sobre este conjunto.

7.2. Espacios de Hilbert con kernel reproductivo en aprendizaje de máquina

Es de conocimiento general que el estudio de funciones kernels definidas positivos es un tema de interés para la comunidad del aprendizaje de máquina como una generalización de un buen cuerpo teórico que ha sido desarrollado para modelos lineales. Una función kernel definida positiva κ , es un camino implícito para representar las muestras de un espacio de entrada \mathcal{X} . Debido a que hay una correspondencia entre κ y los RKHS de funciones en \mathcal{H} , la función kernel puede ser entendida como un camino indirecto para el calculo del producto interno entre elementos de un espacio de Hilbert, los cuales son resultado de elementos mapeados de \mathcal{X} a \mathcal{H} . Por ende, existe una función de mapeo $\varphi : \mathcal{X} \rightarrow \mathcal{H}$ tal que:

$$\kappa(x, x') = \langle \varphi(x), \varphi(x') \rangle_{\mathcal{H}}. \quad (7.4)$$

Con respecto a esto, el espacio \mathcal{H} puede ser visto como un espacio de características y φ es conocida la función de mapeo al espacio de características. Como consecuencia, mediante la realización de operaciones lineales en \mathcal{H} es posible realizar manipulaciones lineales en el espacio de entrada \mathcal{X} ; no obstante, sin necesidad de realizar cualquier calculo explicito en \mathcal{H} (ver figura 7.1).

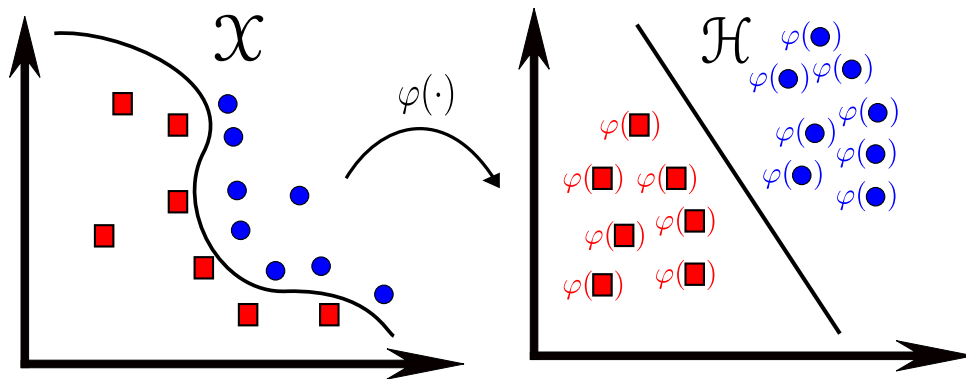


Figura 7.1.: Mapeo basado en función kernel.

Entonces, una importante propiedad asociada con el uso de funciones kernels, definidas positivas, en aprendizaje de máquina, es conocido como el *teorema de representación* [11, 46].

Teorema 7.5. Sea $\Omega : [0, +\infty) \rightarrow \mathbb{R}$ una función monótonamente creciente de forma estricta, \mathcal{X} un conjunto, y $\epsilon : (\mathcal{X} \times \mathbb{R}^2)^N \rightarrow \mathbb{R} \cup \infty$ una función de pérdida arbitraria. Entonces, cada minimizador $f \in \mathcal{H}$ del riesgo regularizado funcional:

$$\epsilon((\mathbf{x}_1, t_1, f(\mathbf{x}_1)), \dots, (\mathbf{x}_N, t_N, f(\mathbf{x}_N))) + \Omega(\|f\|_{\mathcal{H}}^2) \quad (7.5)$$

admite una representación de la forma

$$f(\mathbf{x}) = \sum_{n=1}^N \alpha_n \kappa(\mathbf{x}_n, \mathbf{x}), \quad (7.6)$$

donde cada $t_n \in \mathbb{R}$ es un salida asociada con la entrada $\mathbf{x}_n \in \mathcal{X}$.

Demostración. Sea $S = \text{span}\{\kappa(\mathbf{x}_n, \cdot) : \mathbf{x}_n \in \mathcal{X}, n \in \{1, \dots, N\}\}$ el subespacio extendido \mathcal{H} para las N muestras de entrenamiento. Considere la solución $f \in \mathcal{H}$, esta solución puede ser escrita como: $f = f_S + f_{S^\perp}$, donde $f_S \in S$, $f_{S^\perp} \in S^\perp$, y \perp indica el símbolo ortogonal. Consecuentemente, $f(\mathbf{x}_n) = f_S(\mathbf{x}_n) + f_{S^\perp}(\mathbf{x}_n) = f_S(\mathbf{x}_n) + 0$. Ahora, para el segundo termino del riesgo regularizado funcional:

$$\Omega(\|f\|_{\mathcal{H}}^2) = \Omega(\|f_S\|_{\mathcal{H}}^2 + \|f_{S^\perp}\|_{\mathcal{H}}^2),$$

como Ω es estrictamente monótona creciente, es posible ver que el mínimo será registrado por $\|f_{S^\perp}\|_{\mathcal{H}} = 0$, lo cual implica que $f_{S^\perp} = 0$. □

Con esto en mente, es posible concluir que el teorema de representación básicamente declara que la solución de la función de riesgo regularizado puede ser expresada en términos de las llamadas muestras de entrenamiento $\{\mathbf{x}_n, t_n\}_{n=1}^N$. Por lo tanto, nos permite abordar con problemas que a simple vista parecen ser de dimensión infinita. No obstante, la regularización previene de múltiples mínimos locales, tal propiedad requiere algunas condiciones extras, normalmente, convexidad.

8. Multiplicadores de Lagrange

Los *multiplicadores de Lagrange* son usados para encontrar puntos estacionarios de una función sujeta a varias variables para una o más restricciones. Consideremos el problema de hallar el máximo de la función $f(x_1, x_2)$ sujeto a una restricción relativa a x_1 y x_2 , el cual vamos a escribir de la forma:

$$g(x_1, x_2) = 0. \quad (8.1)$$

Un enfoque sería resolver la restricción de la ecuación (8.1) y entonces expresar x_2 en función de x_1 de la forma $x_2 = h(x_1)$. Este puede ser entonces sustituido dentro de $f(x_1, x_2)$ para obtener una función dependiente únicamente de x_1 , tomando la forma $f(x_1, h(x_1))$. Luego, el máximo con respecto a x_1 puede ser encontrado por el camino usual de la diferenciación, para obtener el valor estacionario x_1^* , y con este valor calcular $x_2^* = h(x_1^*)$.

Uno de los problemas del anterior enfoque es la dificultad que puede presentar encontrar la solución analítica de la ecuación de la restricción que permite a x_2 ser expresado como una función explícita de x_1 . También, este enfoque trata a x_1 y x_2 de forma diferente y así estropea la simetría natural entre estas variables.

Un más elegante, y con frecuencia más simple, enfoque es basado en la introducción de un parámetro λ llamado un *multiplicador de Lagrange*. Motivaremos esta técnica desde una perspectiva geométrica. Consideremos un vector $\mathbf{x} \in \mathbb{R}^D$ con la ecuación de restricción $g(\mathbf{x}) = 0$, la cual representa una superficie de $(D - 1)$ -dimensiones en el espacio de \mathbf{x} como es indicado en la figura 8.1(a).

Primero debemos notar que cualquier punto, sobre la superficie de la restricción, evaluado en el gradiente $\nabla g(\mathbf{x})$ será ortogonal a la superficie generada por $g(\mathbf{x}) = 0$. Para observar esto, considere un punto \mathbf{x} que se encuentra sobre la superficie de restricción, y considere un punto vecino $\mathbf{x} + \epsilon$ que también se encuentra sobre la superficie. Si realizamos una expansión de Taylor alrededor de \mathbf{x} , tenemos

$$g(\mathbf{x} + \epsilon) \simeq g(\mathbf{x}) + \epsilon^\top \nabla g(\mathbf{x}). \quad (8.2)$$

Debido a que ambos \mathbf{x} y $\mathbf{x} + \epsilon$ se encuentran sobre la superficie de restricción, tenemos que $g(\mathbf{x}) = g(\mathbf{x} + \epsilon) = 0$, y por lo tanto $\epsilon^\top \nabla g(\mathbf{x}) \simeq 0$. Si el límite $\|\epsilon\|_2 \rightarrow 0$, entonces $\epsilon^\top \nabla g(\mathbf{x}) = 0$, y porque ϵ es entonces paralelo a la superficie de restricción, vemos que el vector ∇g es normal a esta superficie.

Luego buscamos un punto \mathbf{x}^* sobre la superficie de restricción tal que $f(\mathbf{x}^*)$ sea máximo. Tal punto debe tener la propiedad de que el vector $\nabla f(\mathbf{x}^*)$ sea también ortogonal a la superficie de restricción, $g(\mathbf{x}) = 0$, como se ilustra en la figura 8.1(a), dado que en cualquier otro caso podríamos incrementar el valor de $f(\mathbf{x}^*)$ moviendo una corta distancia sobre la superficie de restricción. De este modo, $\nabla f(\mathbf{x}^*)$ y $\nabla g(\mathbf{x}^*)$ son vectores anti-paralelos, y entonces debe existir un parámetro λ tal que

$$\nabla f(\mathbf{x}^*) + \lambda \nabla g(\mathbf{x}^*) = 0, \quad (8.3)$$

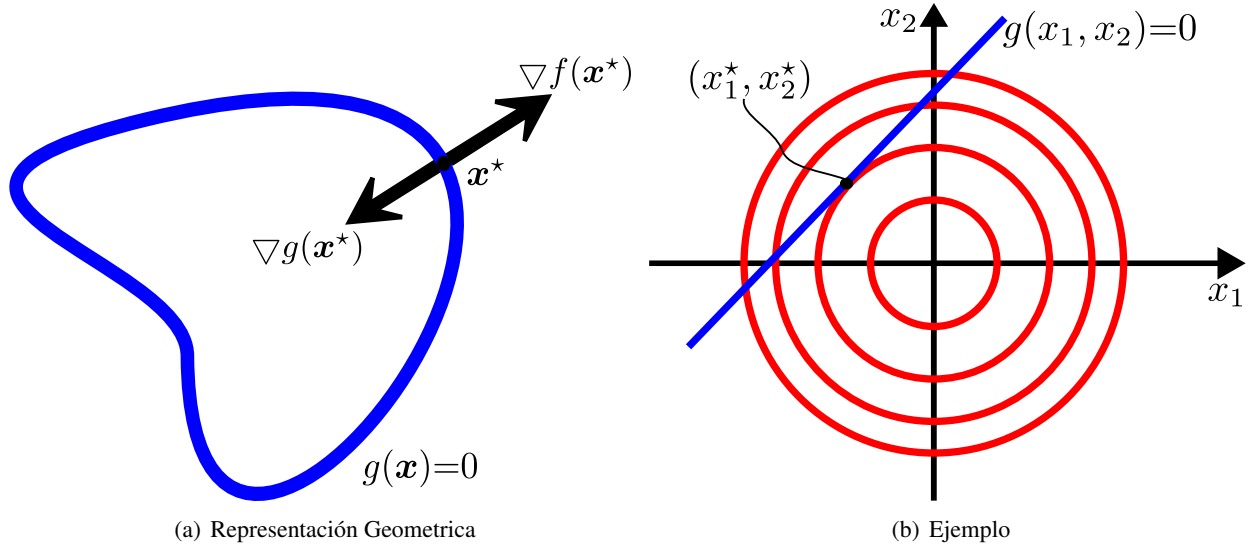


Figura 8.1.: Figura geométrica de la técnica de multiplicadores de Lagrange en el cual buscamos maximizar una función $f(\mathbf{x})$, sujeta a la restricción $g(\mathbf{x})=0$. Si \mathbf{x} es de dimensión D , la restricción corresponde a un subespacio de dimensión $D - 1$, indicada por la curva azul. El problema puede resolverse por optimización de la función Lagrangiana $\mathcal{L}(\mathbf{x}, \lambda)=f(\mathbf{x}) + \lambda g(\mathbf{x})$. Derecha: Un ejemplo sencillo del uso de los multiplicadores de Lagrange en el cual tiene el objetivo de maximizar $f(x_1, x_2)$ sujeto a la restricción $g(x_1, x_2)=0$. Los círculos rojos muestran las curvas de nivel de la función $f(x_1, x_2)$, y la línea azul corresponde a la superficie de restricción $g(x_1, x_2)=0$.

donde $\lambda \neq 0$ es un *multiplicador de Lagrange*. Notemos que λ puede tener cualquier signo.

En este punto, es conveniente introducir la función del *Lagrangiano* definida por:

$$\mathcal{L}(\mathbf{x}, \lambda) \equiv f(\mathbf{x}) + \lambda g(\mathbf{x}). \quad (8.4)$$

La condición de la restricción estacionaria (ecuación (8.3)) es obtenida por medio de la igualdades $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \lambda)=0$ y $\partial \mathcal{L}(\mathbf{x}, \lambda)/\partial \lambda=0$.

De este modo, para encontrar el máximo de la función $f(\mathbf{x})$ sujeta a la restricción $g(\mathbf{x})=0$, definimos la función del Lagrangiano dada en la ecuación (8.4) y entonces encontramos el punto estacionario de $\mathcal{L}(\mathbf{x}, \lambda)$ con respecto a \mathbf{x} y λ . Para un vector $\mathbf{x} \in \mathbb{R}^D$, este otorga $D + 1$ ecuaciones las cuales determinan el punto estacionario \mathbf{x}^* y el valor del multiplicador de Lagrange λ .

Como un simple ejemplo, supongamos que deseamos encontrar el punto estacionario de la función $f(x_1, x_2)=1 - x_1^2 - x_2^2$ sujeto a la restricción $g(x_1, x_2)=x_1 + x_2 - 1=0$, como se ilustra en la figura 8.1(b). La correspondiente función del Lagrangiano esta dada por:

$$\mathcal{L}(\mathbf{x}, \lambda) = 1 - x_1^2 - x_2^2 + \lambda(x_1 + x_2 - 1). \quad (8.5)$$

Las condiciones para que este Lagrangiano sea estacionario con respecto a x_1 , x_2 , y λ están dadas por las siguientes ecuaciones:

$$-2x_1 + \lambda = 0 \quad (8.6)$$

$$-2x_2 + \lambda = 0 \quad (8.7)$$

$$x_1 + x_2 - 1 = 0. \quad (8.8)$$

Solucionando las ecuaciones obtenemos el punto estacionario como $[x_1^*, x_2^*] = [\frac{1}{2}, \frac{1}{2}]$, y el correspondiente valor para el multiplicador de Lagrange es $\lambda=1$.

Hasta el momento, hemos considerado un problema de minimización sujeta a una *restricción de igualdad* de la forma $g(x)=0$. Ahora consideraremos el problema de maximizar $f(x)$ sujeta a una *restricción de desigualdad* de la forma $g(x) \geq 0$, como se ilustra en la figura 8.2.

Ahora existen dos tipos de soluciones posibles, acorde a la restricción estacionaria, que son los puntos que se encuentran dentro de la región $g(x) > 0$, en tal caso la restricción es *inactiva*, o los que se encuentran sobre la restricción $g(x)=0$, en este caso se dice que la restricción es *activa*. En el primer caso, la función $g(x)$ no juega ningún papel y la condición estacionaria es simplemente $\nabla f(x)=0$. Este de nuevo corresponde al punto estacionario de la función de Lagrange definida en la ecuación (8.4), pero con $\lambda=0$. El otro caso, donde la solución está sobre el límite, es análogo a la restricción de igualdad discutido previamente y corresponde al punto estacionario de la función de Lagrange con $\lambda \neq 0$. Sin embargo, el signo del multiplicador de Lagrange es crucial, debido a que la función $f(x)$ sólo puede ser máxima si su gradiente es orientado lejos de la región $g(x) > 0$, como se ilustra en la figura 8.2. Por lo tanto, tenemos $\nabla f(x) = -\lambda \nabla g(x)$ para algunos valores de $\lambda > 0$.

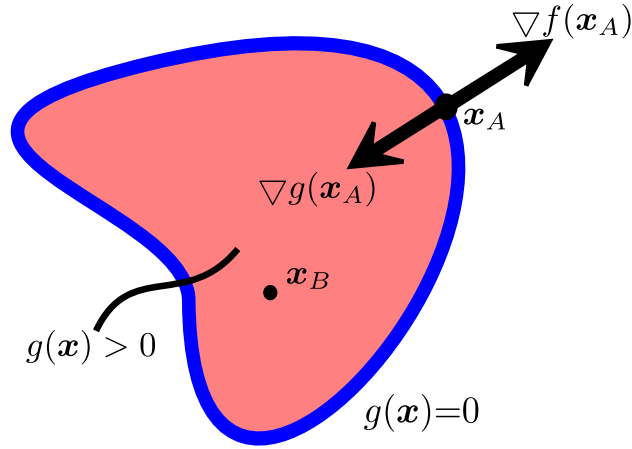


Figura 8.2.: Representación de el problema de minimizar $f(x)$ sujeta a la restricción de desigualdad $g(x) \geq 0$.

Para cualquiera de estos dos casos, el producto $\lambda g(x)=0$. Así, la solución para el problema de minimización de $f(x)$ sujeta a la restricción $g(x) \geq 0$ es obtenida por la optimización de la función de Lagrange (ecuación (8.4)) con respecto a x y λ sujeta a las condiciones

$$g(x) \geq 0 \quad (8.9)$$

$$\lambda \geq 0 \quad (8.10)$$

$$\lambda g(x) = 0. \quad (8.11)$$

Estas condiciones son conocidas como las *condiciones de Karush-Kuhn-Tucker* (KKT) [47, 48].

Notemos que si deseamos minimizar la función $f(\mathbf{x})$ sujeto a la restricción de desigualdad $g(\mathbf{x}) \geq 0$, sólo es necesario minimizar la función del Lagrangiano $\mathcal{L}(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda g(\mathbf{x})$ con respecto a \mathbf{x} , y de nuevo sujeto a que $\lambda \geq 0$.

Finalmente, para extender la técnica de multiplicadores de Lagrange al caso de múltiples restricciones de igualdad y desigualdad, suponemos que deseamos maximizar la función $f(\mathbf{x})$ sujeto a $g_j(\mathbf{x}) = 0$ para $j \in \{1, 2, \dots, J\}$, y $h_k(\mathbf{x}) \geq 0$ para $k \in \{1, 2, \dots, K\}$. Introducimos entonces los multiplicadores de Lagrange $\{\lambda_j\}_{j=1}^J$ y $\{\mu_k\}_{k=1}^K$, y entonces optimizamos la función del Lagrangiano dada por

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{j=1}^J \lambda_j g_j(\mathbf{x}) + \sum_{k=1}^K \mu_k h_k(\mathbf{x}), \quad (8.12)$$

sujeto a que $\mu_k \geq 0$ y $\mu_k h_k(\mathbf{x}) = 0$ para $k \in \{1, \dots, K\}$.

9. Preliminares Clasificación

9.1. Aprendizaje supervisado, clasificación desbalanceadas

Para nosotros es fácil reconocer un rostro, entender palabras habladas, leer caracteres realizados a mano alzada y decidir sobre una naranja o una mandarina por su olor; esto, ha sido indispensable para nuestra supervivencia en los pasados diez mil años. El reconocimiento de patrones se define por [49] como: “*el acto de tomar los datos sin procesar y realizar una acción basada en categoría del patrón*”. Debido, al avance tecnológico hemos desarrollado sistemas neuronales y cognitivos sofisticados para realizar esta tarea a gran escala.

Dicho lo anterior, dentro del reconocimiento de patrones se encuentra la tarea de clasificación, la cual consiste en asignar una categoría o clase a un elemento, o muestra, a partir de características de los datos. Un ejemplo de clasificación binaria, es el problema de un granjero para separar por medio de una cerca a sus caballos de sus vacas con el fin de tenerlos en dos zonas de pastoreo diferentes. El inconveniente es determinar cuál es la mejor dirección de la cerca para esta tarea. La figura 9.1 ejemplifica la situación mencionada; podemos observar que existen varias opciones para separar ambas especies, una de ellas es por medio de la cerca azul y otra por la roja, esto nos deja la pregunta ¿cuál es la mejor opción?. Y así determinar a partir de la zona si un animal es de una especie o de la otra.

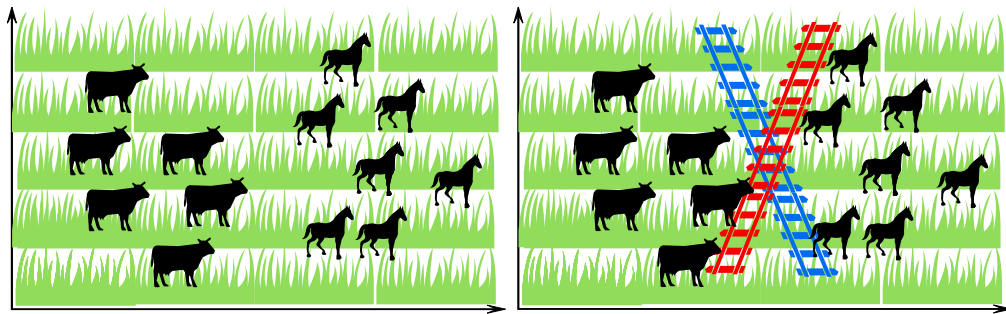


Figura 9.1.: Ejemplo de clasificación

Para realizar esta tarea desde un punto de vista estadístico y/o probabilístico se hace uso del aprendizaje de máquina, el cual definimos como un conjunto de métodos que pueden detectar patrones en datos de forma automática, y que luego usa estos patrones detectados para predecir un dato futuro (nuevo), o realizar otro tipo de toma de decisiones bajo incertidumbre; por ejemplo, decidir si un dato pertenece a alguna clase en específico [50]. En este trabajo nos concentramos en las famosas Máquinas de Vectores de Soporte (SVM, por sus siglas en inglés) propuestas por [51]. Este método de clasificación y regresión es famoso por su capacidad de generalización y su uso en diferentes tareas del mundo real como identificación de

partículas, categorización de texto, bio-informática, y aplicaciones financieras [52–54]. Este algoritmo de aprendizaje en su etapa de entrenamiento debe resolver un problema de optimización cuadrático (*quadratic programming problem*–(QPP)), además su complejidad computacional es de $\mathcal{O}(N^3)$, siendo N el número de datos a utilizar en esta etapa.

9.2. Máquinas de vectores de Soporte (SVM)

Sea $\{\mathbf{X} \in \mathbb{R}^{P \times N}, \mathbf{t} \in \{-1, +1\}^N\}$ una base de datos con N muestras de entrada $\mathbf{x}_n \in \mathbb{R}^P$ y salidas etiquetadas $t_n \in \{-1, +1\}$, con $n \in \{1, 2, \dots, N\}$, para clasificación binaria. La SVM busca encontrar un hiperplano de la forma:

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b = 0, \quad (9.1)$$

donde $\mathbf{w} \in \mathbb{R}^P$ es un vector normal ponderado del hiperplano y $b \in \mathbb{R}$ es un termino de sesgo. Además, la SVM escoge la frontera de decisión (hiperplano) que maximizá la margen, la cual es definida como la mínima distancia perpendicular entre el hiperplano y cualquiera de las muestras. Pero, primero debemos definir la expresión para la distancia perpendicular de una muestra \mathbf{x} al hiperplano. Podemos demostrar que el valor de $f(\mathbf{x})$ da una medida con signo de la distancia perpendicular r del punto \mathbf{x} desde la superficie de decisión. Para observar esto, consideremos el escenario propuesto por la figura 9.2.

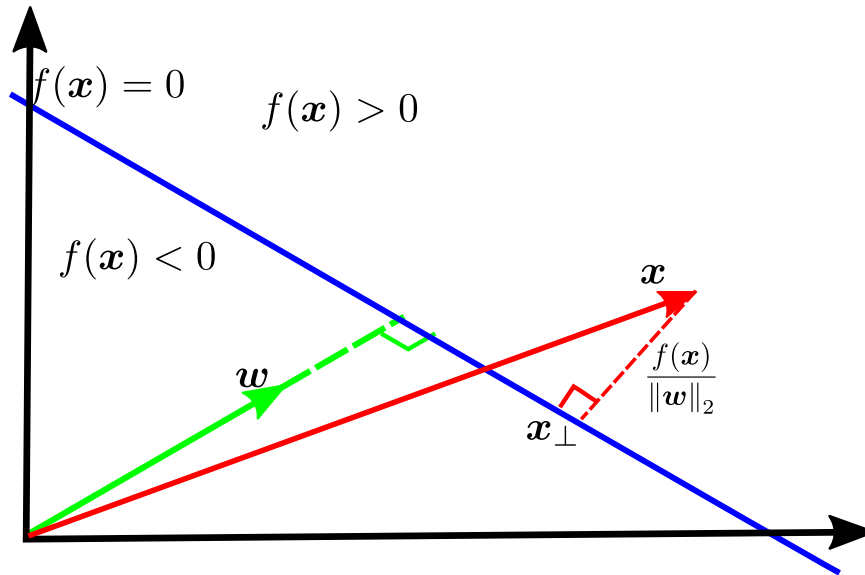


Figura 9.2.: Ilustración geométrica de la distancia ortogonal de una muestra \mathbf{x} , representada en color rojo y la superficie de decisión, vista en azul; su vector perpendicular \mathbf{w} es ilustrado en verde.

También consideremos un punto arbitrario \mathbf{x} y sea \mathbf{x}_\perp su proyección ortogonal sobre el hiperplano $f(\mathbf{x})=0$, entonces podemos decir que

$$\mathbf{x} = \mathbf{x}_\perp + r \frac{\mathbf{w}}{\|\mathbf{w}\|_2}, \quad (9.2)$$

donde $\|\cdot\|_2$ es el operador de norma L_2 .

Ahora multiplicando ambos lados por el término \mathbf{w}^\top y luego sumamos b a ambos lados, y usando las expresión $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ y $f(\mathbf{x}_\perp) = \mathbf{w}^\top \mathbf{x}_\perp + b = 0$, obtenemos que

$$r = \frac{f(\mathbf{x})}{\|\mathbf{w}\|_2}. \quad (9.3)$$

Ahora, esta expresión puede llegar a poseer signo negativo dependiendo de la ubicación del punto \mathbf{x} ; por esto, hacemos uso del operador valor absoluto $|\cdot|$ y la expresión para calcular la distancia perpendicular de una muestra al hiperplano queda de la forma

$$\frac{|f(\mathbf{x})|}{\|\mathbf{w}\|_2} = \frac{t_n f(\mathbf{x})}{\|\mathbf{w}\|_2} = \frac{t_n (\mathbf{w}^\top \mathbf{x} + b)}{\|\mathbf{w}\|_2}. \quad (9.4)$$

Siguiendo con la formulación en la SVM, la figura 9.3(a) muestra un ejemplo gráfico de la mismas e ilustra la definición de margen. También, observamos muestras pertenecientes a dos clases, una clase representada con el color rojo y otra por el azul. Podemos ver en esta figura, que la margen la define la muestra roja más cercana, ortogonalmente, al hiperplano; adicionalmente, este punto de alguna manera definirá el hiperplano y por esto vendría hacer parte del subconjunto de muestras llamado, *vectores de soporte*—(SV, por sus siglas en inglés).

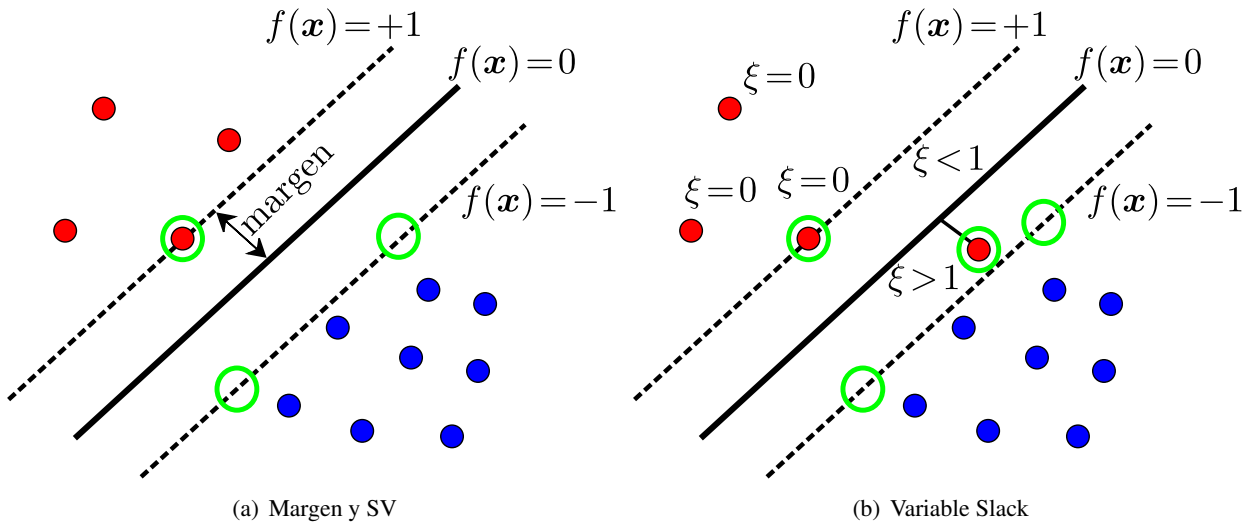


Figura 9.3.: La SVM busca maximizar la margen, la cual se define como la mínima distancia perpendicular del hiperplano a cualquiera de las muestras disponibles. (a) muestra un ejemplo gráfico sobre la definición del la *magen*. (b) es un ejemplo donde existe traslape entre las clases y se permite un hiperplano que clasifique muestras de una clase que no se encuentren en zona correspondiente, para ello se hace uso de las variables auxiliares ξ_n . En ambas figuras, las muestras encerradas en un círculo verde son consideradas vectores de soporte y son las que determinan el hiperplano.

Por otra parte, cuando las muestras de ambas clases se traslapan definir una frontera de decisión lineal es una tarea casi imposible. Por lo anterior, las SVM consideran variables auxiliares $\xi_n = |t_n - f(\mathbf{x}_n)|$, para $n \in \{1, 2, \dots, N\}$. Las muestras correspondientes a $\xi_n = 0$ estarían bien clasificadas y a una distancia mayor que la margen; por otra parte, las muestras donde $\xi_n = 1$ indicaría que están sobre la frontera de decisión; y por último, las muestras con $\xi_n > 1$, quieren decir que se encuentran mal clasificadas por la frontera de decisión. Estas variables auxiliares dan a este modelo de una tolerancia ante la mala clasificación, de algunas muestras difíciles, como lo son por ejemplo las que se traslapan con la otra clase, y así generar una frontera de decisión más generalizante. Por este motivo, el conjunto de muestras que definen la margen con $\xi_n > 1$

son considerados vectores de soporte.

Asimismo, nos interesa que todos las muestras queden correctamente clasificadas, para que $t_n f(\mathbf{x}_n) > 0 \forall n=1,2,\dots,N$. La margen esta dada como la distancia perpendicular al punto más cercano \mathbf{x}_n del conjunto de datos, y deseamos optimizar los parámetros \mathbf{w} y b para maximizar esta distancia, Así, la solución para la máxima margen se encuentra resolviendo,

$$\mathbf{w}^*, b^* = \arg \max_{\mathbf{w}, b} \left\{ \frac{1}{\|\mathbf{w}\|_2} \min_n \left[t_n (\mathbf{w}^\top \mathbf{x}_n + b) \right] \right\} \quad (9.5)$$

donde el factor $1/\|\mathbf{w}\|_2$ se factoriza del problema de optimización sobre n , dado que no depende de n . La solución directa de este problema de optimización es compleja; por ende, se cambia el problema de optimización de la margen a uno más sencillo de resolver. Para lo anterior, debemos de notar que si re-escalamos los parámetros del hiperplano, $\mathbf{w} \rightarrow \alpha \mathbf{w}$ y $b \rightarrow \alpha b$, entonces la distancia de cualquier muestra \mathbf{x}_n a la frontera de decisión, $t_n f(\mathbf{x}_n)/\|\mathbf{w}\|_2$, no cambia y todas las muestras satisfarán la restricción

$$t_n (\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 - \xi_n, \quad \forall n \in \{1, 2, \dots, N\}. \quad (9.6)$$

Lo anterior, es conocido como representación canónica del hiperplano de decisión. En el caso de las muestras que mantienen la igualdad, se dice que las restricciones están activas, mientras que para las demás muestras se dice que están inactivas. Por la definición de la margen, siempre habrá al menos una muestra activa. El problema de optimización simplemente requerirá que maximicemos $\|\mathbf{w}\|_2^{-1}$, lo cual es equivalente a minimizar $\|\mathbf{w}\|_2^2$, entonces tenemos que resolver el siguiente problema de optimización:

$$\begin{aligned} \tilde{\mathbf{w}}, \tilde{b} = \arg \min_{\mathbf{w}, b, \xi} & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{n=1}^N \xi_n \\ \text{s.t. } & t_n (\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 - \xi_n, \\ & \xi_n \geq 0 \end{aligned} \quad (9.7)$$

$\forall n \in \{1, 2, \dots, N\}$; donde $C \in \mathbb{R}^+$ es el parámetro que regulariza la penalización de las variables auxiliares y la margen. Debido a que cualquier punto mal clasificado tiene $\xi_n > 1$, además es importante recalcar que al minimizar el termino $\|\mathbf{w}\|_2^2$ junto con la restricción expresada en la ecuación (9.6), se maximiza de una manera menos compleja la margen del hiperplano; por otra parte, el minimizar el termino $\sum_n \xi_n$ ayuda a mitigar los escenarios de traslape entre clases y así reducir en lo posible la mala clasificación de las muestras de entrenamiento. El lagrangiano (ver capítulo 8) del problema de optimización expresado en la ecuación (9.7) esta dado por:

$$\mathcal{L}(\mathbf{w}, b, \xi) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} + C \sum_{n=1}^N \xi_n - \sum_{n=1}^N \alpha_n (t_n f(\mathbf{x}_n) - 1 + \xi_n) - \sum_{n=1}^N \mu_n \xi_n, \quad (9.8)$$

donde $\{\alpha_n \geq 0\}$ y $\{\mu_n \geq 0\}$ son multiplicadores de Lagrange, y $\xi \in \mathbb{R}^N$ es el vector de variables auxiliares con elementos ξ_n . Su correspondiente conjunto de condiciones *Karush-Kuhn-Tucker* (KKT) (ver capítulo 8) son expresadas por:

$$\alpha_n \geq 0 \quad (9.9)$$

$$t_n f(\mathbf{x}_n) - 1 + \xi_n \geq 0 \quad (9.10)$$

$$\alpha_n (t_n f(\mathbf{x}_n) - 1 + \xi_n) = 0 \quad (9.11)$$

$$\mu_n \geq 0 \quad (9.12)$$

$$\xi_n \geq 0 \quad (9.13)$$

$$\mu_n \xi_n = 0, \quad (9.14)$$

$\forall n \in \{1, 2, \dots, N\}$.

Ahora optimizamos \mathbf{w} , b , y $\{\xi_n\}_{n=1}^N$ haciendo uso de la definición para $f(\mathbf{x})$ dada en la ecuación (9.1), con lo que obtenemos:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0} \Rightarrow \mathbf{w} = \sum_{n=1}^N \alpha_n t_n \mathbf{x}_n \quad (9.15)$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow \sum_{n=1}^N \alpha_n t_n = 0 \quad (9.16)$$

$$\frac{\partial \mathcal{L}}{\partial \xi_n} = 0 \Rightarrow \alpha_n = C - \mu_n. \quad (9.17)$$

Usando estos resultados con el fin de eliminar \mathbf{w} , b y $\{\xi_n\}_{n=1}^N$ de $\mathcal{L}(\mathbf{w}, b, \xi)$, obtenemos la *representación dual* del problema de optimización (ecuación (9.7)), expresada por:

$$\begin{aligned} \tilde{\alpha} = \arg \max_{\alpha} \quad & \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m \mathbf{x}_n^\top \mathbf{x}_m \\ \text{s.t. } & 0 \leq \alpha_n \leq C \\ & \sum_{n=1}^N \alpha_n t_n = 0, \end{aligned} \quad (9.18)$$

$\forall n \in \{1, 2, \dots, N\}$. La primera restricción esta formulada a partir de su naturaleza como multiplicadores de Lagrange, $\{\alpha_n \geq 0\}_{n=1}^N$, y para asegurar que $\{\mu_n \geq 0\}_{n=1}^N$ es evidente que $\{\alpha_n \leq C\}_{n=1}^N$ (ver ecuación (9.17)). Y la segunda condición es acorde a la derivada de $\mathcal{L}(\mathbf{w}, b, \xi)$ con respecto al parámetro b (ecuación (9.16)). Es importante recalcar que el problema de optimización dual es un *problema de programación cuadrático-QPP* de una sola variable, α^* .

Después de lo anterior, para clasificar nuevas muestras \mathbf{x}_* usando el modelo entrenado evaluamos la función signo sobre $f(\mathbf{x}_*)$ definido en la ecuación (9.1), $\text{sign}(f(\mathbf{x}_*))$. El expresión del hiperplano, $f(\mathbf{x})$, puede ser expresado en termino de los $\{\tilde{\alpha}_n\}_{n=1}^N$ (donde $\tilde{\alpha}_n$ es el n -ésimo elemento de $\tilde{\alpha}$) y del producto interior entre la nueva muestra y las muestras de entrenamiento, sustituyendo \mathbf{w} usando la ecuación (9.15), para obtener:

$$f(\mathbf{x}_*) = \sum_{n=1}^N \tilde{\alpha}_n t_n \langle \mathbf{x}_*, \mathbf{x}_n \rangle + b \quad (9.19)$$

En el capítulo 8, mostramos que un problema de optimización restringida de esta forma satisface las condiciones *Karush-Huhn-Tucker*-(KKT), y se requiere que las propiedades

$$\alpha_n \geq 0 \quad (9.20)$$

$$t_n f(\mathbf{x}_n) - 1 + \xi_n \geq 0 \quad (9.21)$$

$$\alpha_n (t_n f(\mathbf{x}_n) - 1 + \xi_n) = 0, \quad (9.22)$$

se cumplan para todo n . Asimismo, para todas las muestras puede suceder que $\alpha_n=0$ ó $t_n f(\mathbf{x}_n)=1$. Dicho lo anterior, cualquier muestra para la cual $\alpha_n = 0$, no contribuye a la sumatoria en la ecuación (9.19) y por ende tiene un papel en la predicción de nuevas muestras. Por otra parte, las muestras restantes, $\alpha_n \neq 0$, son llamadas *vectores de soporte* (SV), porque corresponden a las muestras que se encuentran sobre la margen generada y que se encuentran del lado equivocado del hiperplano (esto se ilustra en la figura 9.3); todas estas muestras satisfacen la condición $t_n f(\mathbf{x}_n)=1 - \xi_n$. La propiedad anterior es central para la aplicabilidad de la SVM. Una vez el modelo es entrenado, una proporción significativa de las muestras pueden ser descartada y sólo los vectores de soporte son conservados.

Habiendo resultado el QPP y encontrado los valores de $\tilde{\alpha}$, podemos entonces determinar el valor del parámetro de sesgo b notando que cualquier vector de soporte \mathbf{x}_n satisface $t_n f(\mathbf{x}_n)=1$. Usando la ecuación (9.19) obtenemos:

$$t_n \left(\sum_{m \in S} \tilde{\alpha}_m t_m \langle \mathbf{x}_n, \mathbf{x}_m \rangle + b \right) = 1, \quad (9.23)$$

donde S es el conjunto de índices de los vectores de soporte. Aunque podemos solucionar esta ecuación para b usando un vector de soporte \mathbf{x}_n escogido arbitrariamente, una solución más estable numéricamente es obtenida multiplicando primero por t_n , y como $t_n^2=1$, entonces despejamos b de la expresión; por último promediamos los resultados sobre todas las muestras que son vectores de soporte, obteniendo:

$$b = \frac{1}{N_S} \sum_{n \in S} \left(t_n - \sum_{m \in S} \tilde{\alpha}_m t_m \langle \mathbf{x}_n, \mathbf{x}_m \rangle \right), \quad (9.24)$$

donde N_S es el número de vectores de soporte total.

9.2.1. Máquinas de Vectores de Soporte en RKHS

Para extender las SVM a problema de clasificación no-lineales, se busca representar los datos en un espacio de alta dimensionalidad por medio de una función de mapeo no-lineal $\varphi(\cdot) : \mathbb{R}^P \rightarrow \mathbb{R}^D$, con $D \rightarrow \infty$. La idea es la misma a la del clasificador explicado anteriormente, generar un hiperplano de separación pero en un espacio de características de alta dimensionalidad (\mathbb{R}^D) en este orden de ideas, el hiperplano toma la forma:

$$f(\mathbf{x}) = \mathbf{w}^\top \varphi(\mathbf{x}) + b = 0. \quad (9.25)$$

Siguiendo el mismo procedimiento del apartado 9.2 y usando el mismo problema de optimización primario :

$$\begin{aligned} \tilde{\mathbf{w}}, \tilde{b} = \arg \min_{\mathbf{w}, b, \xi} & \frac{1}{2} \|\mathbf{w}\|_2^2 + C \sum_{n=1}^N \xi_n \\ \text{s.t. } & t_n (\mathbf{w}^\top \varphi(\mathbf{x}_n) + b) \geq 1 - \xi_n, \\ & \xi_n \geq 0, \end{aligned} \quad (9.26)$$

para todo n , donde $\varphi(\mathbf{x}_n) \in \mathbb{R}^D$ es la función de mapeo no-lineal evaluada en la n -ésima muestra, $\mathbf{w} \in \mathbb{R}^D$ es el vector de pesos, ortogonal al hiperplano, y $b \in \mathbb{R}$ es el término de sesgo. Siguiendo el mismo procedimiento de la sección anterior y cambiando \mathbf{x}_n por $\varphi(\mathbf{x}_n)$, y \mathbf{x}_* por $\varphi(\mathbf{x}_*)$, llegamos al problema dual de la forma:

$$\begin{aligned} \tilde{\alpha} = \arg \max_{\alpha} \quad & \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m \kappa(\mathbf{x}_n, \mathbf{x}_m), \\ \text{s.t. } \quad & 0 \leq \alpha_n \leq C, \\ & \sum_{n=1}^N \alpha_n t_n = 0, \end{aligned} \quad (9.27)$$

donde $\kappa : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$ es una función kernel (ver apartado 7.2), con $\langle \varphi(\mathbf{x}_n), \varphi(\mathbf{x}_m) \rangle = \kappa(\mathbf{x}_n, \mathbf{x}_m)$. De nuevo, el problema de optimización dual toma la forma de un QPP en función de solamente α .

En este orden de ideas, el hiperplano definido en la ecuación (9.25) en función de los $\{\tilde{\alpha}_n\}_{n=1}^N$ para predecir una nueva muestra \mathbf{x}_* , toma la siguiente forma:

$$f(\mathbf{x}_*) = \sum_{n=1}^N \tilde{\alpha}_n t_n \kappa(\mathbf{x}_*, \mathbf{x}_n) + b, \quad (9.28)$$

y la expresión para calcular el termino de sesgo b queda como:

$$b = \frac{1}{N_S} \sum_{n \in S} \left(t_n - \sum_{m \in S} \tilde{\alpha}_m t_m \kappa(\mathbf{x}_n, \mathbf{x}_m) \right). \quad (9.29)$$

Dicho todo lo anterior, es importante recordar que la solución de un QPP de P variables, en general tiene una complejidad computacional de $\mathcal{O}(P^3)$. Si vamos a la formulación dual a la que hemos llegado a partir del problema de optimización original, el cual involucra problemas de optimización (ver ecuaciones (9.7) y (9.26)) sobre variables de dimensión P , con N variables. Para un conjunto muestras fijas de dimensión P más pequeño que el número de muestras N , el procedimiento para llegar a una formulación dual parecería que posee desventajas. Sin embargo, permite al modelo ser reformulado usando kernels, por lo que el clasificador con margen máxima puede ser aplicado eficientemente en un espacio de características cuya dimensionalidad, $D \rightarrow \infty$, excede el número de muestras. La formulación del kernel también deja claro el rol de la restricción de la función kernel, tiene que ser definida positiva, ya que esta garantiza que el problema de optimización dual este bien definido.

9.2.2. Máquinas de Vectores de Soporte ponderadas (WSVM)

Debido que la formulación de la SVM no considera si la base de datos con la cual se desea entrenar es balanceada o no, puede ocurrir algún sesgo hacia alguna de las clases [6]. Además, dentro de varias aplicaciones del mundo real, una de las clases tiene una mayor penalización si el clasificador llega a equivocarse en ella en comparación a la otra clase. Por estos motivos, Osuna et al. en [55] proponen el clasificador *Weighted Support Vectors Machine* (WSVM), el cual regulariza las variables auxiliares, $\{\xi_n\}$, dependiendo de la etiqueta, t_n . Dicho esto, el problema primal toma la forma:

$$\begin{aligned}
 \tilde{\mathbf{w}}, \tilde{b} = \arg \min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|_2^2 + C_+ \sum_{t_n=+1} \xi_n + C_- \sum_{t_n=-1} \xi_n \\
 \text{s.t. } t_n (\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 - \xi_n, \\
 \xi_n \geq 0,
 \end{aligned} \tag{9.30}$$

donde $C_+, C_- \in \mathbb{R}^+$ los parámetros de regularización de las variables auxiliares para las muestras de la clase $+1$ y -1 , respectivamente. Siguiendo el mismo procedimiento del apartado 9.2 llegamos al problema dual:

$$\begin{aligned}
 \tilde{\alpha} = \arg \max_{\alpha} \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m \mathbf{x}_n^\top \mathbf{x}_m \\
 \text{s.t. } 0 \leq \alpha_n \leq C_+ \quad \forall n \mid t_n = +1 \\
 0 \leq \alpha_{n'} \leq C_- \quad \forall n' \mid t_{n'} = -1 \\
 \sum_{n=1}^N \alpha_n t_n = 0;
 \end{aligned} \tag{9.31}$$

si suponemos un hiperplano en el espacio original de la forma de la ecuación (9.25), o si se supone un hiperplano en un espacio de características de alta dimensionalidad dado por la ecuación (9.25), el problema dual toma la forma:

$$\begin{aligned}
 \tilde{\alpha} = \arg \max_{\alpha} \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m \kappa(\mathbf{x}_n, \mathbf{x}_m) \\
 \text{s.t. } 0 \leq \alpha_n \leq C_+ \quad \forall n \mid t_n = +1 \\
 0 \leq \alpha_{n'} \leq C_- \quad \forall n' \mid t_{n'} = -1 \\
 \sum_{n=1}^N \alpha_n t_n = 0.
 \end{aligned} \tag{9.32}$$

La ecuación de decisión esta dada por la ecuación (9.19) o por la ecuación (9.28), dependiendo si se usa la SVM en el espacio original o en un espacio de características de alta dimensionalidad.

9.3. Máquinas de Vectores de Soporte Gemelas (TWSVM)

Desde su creación, las SVM han demostrado tener una gran capacidad de generalización y uso en diferentes aplicaciones de varias áreas del conocimiento. Aunque la investigación sobre las SVM ha logrado logros notables, todavía posee varias deficiencias en un estudio más profundo. Por ejemplo, los problemas que incluyen la relación entre la teoría de aprendizaje estadístico y otros sistemas teóricos, el procesamiento de grandes cantidades de datos y la selección de los parámetros. Especialmente, con el rápido desarrollo de la internet y de los sistemas de información, la alta dimensionalidad, la distribución y los datos dinámicos complejos son rápidamente generados.

Con el fin de reducir la complejidad computacional de la SVM, en 2006 Mangasarian y Wild Edward [56] propusieron un clasificador no-lineal, basados en la SVM, nombrada *generalized eigenvalue proximal support vector machine*-(GEPSVM). La idea principal del GEPSVM es buscar dos planos no-paralelos,

exigiendo que las muestras de cada clase estén lo más cerca a su correspondiente plano y a su vez lo más lejos posible de la otra clase. Aunque, la GEPSVM tiene un bajo tiempo de aprendizaje, su *acierto de clasificación* es bajo. En 2007, Jayadeva y Suresh [57] propusieron un nuevo método de aprendizaje de máquina llamado *twin support vector machine*—(TWSVM)—(máquinas de vectores de soporte gemelas, en español) para clasificación binaria con el espíritu del GEPSVM. El TWSVM busca generar dos planos no-paralelos, tal que cada hiperplano encierra a una de las dos clases, lo más lejos posible de la otra.

9.3.1. TWSVM lineal

Ahora, definimos las matrices $\mathbf{X}_+ \in \mathbb{R}^{P \times N_+} = [\mathbf{x}_i \mid t_i = +1]_{i=1}^{N_+}$ y $\mathbf{X}_- \in \mathbb{R}^{P \times N_-} = [\mathbf{x}_j \mid t_j = -1]_{j=1}^{N_-}$, donde $N = N_+ + N_-$, como las matrices de las clases $+1$ y -1 , respectivamente. Observemos que hacemos uso de los índices $+$ y $-$ para hacer referencia a las clases $+1$ y -1 , respectivamente. Como se dijo anteriormente, el TWSVM busca generar dos hiperplanos no-paralelos de la forma:

$$f_+(x) = \mathbf{w}_+ x + b_+ = 0 \quad \text{y} \quad f_-(x) = \mathbf{w}_- x + b_- = 0, \quad (9.33)$$

donde $\mathbf{w}_+, \mathbf{w}_- \in \mathbb{R}^P$ son los vectores normales a sus respectivos hiperplanos y b_+, b_- son los términos de sesgo. Para encontrar los parámetros de ambos hiperplanos, dados en la ecuación (9.33), con el fin de minimizar la distancia de estos a las muestras correspondientes a su clase y que a su vez se encuentren lo más lejos posible de la otra clase, se busca resolver los siguientes problemas de optimización:

$$\begin{aligned} \tilde{\mathbf{w}}_+, \tilde{b}_+, \tilde{\boldsymbol{\xi}}_- &= \arg \min_{\mathbf{w}_+, b_+, \boldsymbol{\xi}_-} \frac{1}{2} \|\mathbf{X}_+^\top \mathbf{w}_+ + \mathbf{1}_+ b_+\|_2^2 + c_{-,1} \mathbf{1}_-^\top \boldsymbol{\xi}_- \\ \text{s.t.} \quad &-1 \left(\mathbf{X}_-^\top \mathbf{w}_+ + \mathbf{1}_- b_+ \right) + \boldsymbol{\xi}_- \geq \mathbf{1}_- \\ &\boldsymbol{\xi}_- \geq 0, \end{aligned} \quad (9.34)$$

$$\begin{aligned} \tilde{\mathbf{w}}_-, \tilde{b}_-, \tilde{\boldsymbol{\xi}}_+ &= \arg \min_{\mathbf{w}_-, b_-, \boldsymbol{\xi}_+} \frac{1}{2} \|\mathbf{X}_-^\top \mathbf{w}_- + \mathbf{1}_- b_-\|_2^2 + c_{+,1} \mathbf{1}_+^\top \boldsymbol{\xi}_+ \\ \text{s.t.} \quad &+1 \left(\mathbf{X}_+^\top \mathbf{w}_- + \mathbf{1}_+ b_- \right) + \boldsymbol{\xi}_+ \geq \mathbf{1}_+ \\ &\boldsymbol{\xi}_+ \geq 0, \end{aligned} \quad (9.35)$$

donde $\boldsymbol{\xi}_+ \in \mathbb{R}^{N_+}$, $\boldsymbol{\xi}_- \in \mathbb{R}^{N_-}$ son los vectores de variables auxiliares de las muestras de la clase $+1$ y -1 respectivamente, $c_{+,1}, c_{-,1} \in \mathbb{R}^+$ son los parámetros de regularización de las variables auxiliares; y, $\mathbf{1}_+ \in \mathbb{R}^{N_+}$ y $\mathbf{1}_- \in \mathbb{R}^{N_-}$ son vectores de unos.

El primer termino de las funciones objetivo, de las ecuaciones (9.34) y (9.35), es la suma al cuadrado del score de las muestras de la clase correspondiente al hiperplano, expresada por los hiperplanos definidos en la ecuación (9.33), la cual es proporcional a la suma al cuadrado de la distancia perpendicular del hiperplano a las muestras (recordar que la distancia perpendicular de un hiperplano a una muestra esta dado por $|t_n f_\ell(x)| / \|\mathbf{w}_\ell\|_2$). La restricciones requieren que el hiperplano tenga mínimo una distancia de 1 a las muestras de la otra clase; además, un conjunto de variables auxiliares son usadas para contrarrestar el inconveniente de que los hiperplanos no puedan tener una distancia de mínimo 1 con las muestras de la otra clase, es decir, con el fin de contrarrestar traslape entre las clases. El segundo termino es la suma de las variables auxiliares, con la meta de minimizar la mala clasificación. Esto se ilustra en la figura 9.4.

De ahora en adelante, con el fin de agilizar el proceso, usaremos las variables índice $\ell \in \{+, -\}$, $\ell' = -\ell$ y ζ_ℓ , donde $\zeta_\ell = +1$ si $\ell = +$ y -1 en caso contrario.

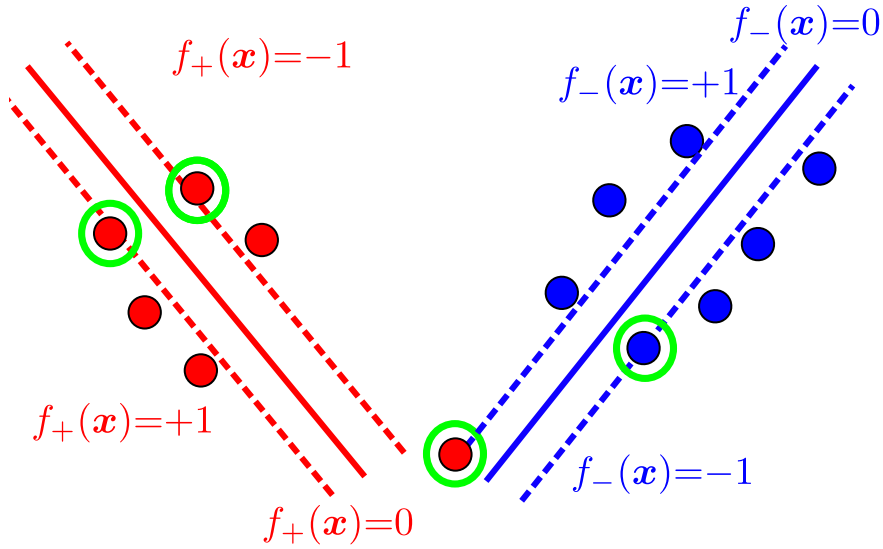


Figura 9.4.: La TWSVM genera dos hiperplanos no-paralelos, donde cada uno está dedicado a una clase. Cada hiperplano debe estar lo más cerca a su clase correspondiente y lo más alejado posible de las muestras de la otra clase. Las muestras e hiperplanos de color rojo corresponden a la clase +1 y los de color azul a la clase -1. Las muestras que se encuentran cerradas en las circunferencias verdes, son las muestras que definen los márgenes de los hiperplanos.

El Lagrangiano de los problemas primales expresados en las ecuaciones (9.34) y (9.35) están dados por:

$$\begin{aligned} \mathcal{L}_\ell(\mathbf{w}_\ell, b_\ell, \boldsymbol{\xi}_{\ell'}) = & \frac{1}{2} (\mathbf{X}_\ell^\top \mathbf{w}_\ell + \mathbf{1}_\ell b_\ell)^\top (\mathbf{X}_\ell^\top \mathbf{w}_\ell + \mathbf{1}_\ell b_\ell) + c_{\ell,1} \mathbf{1}_{\ell'} \cdot \dots \\ & \boldsymbol{\alpha}_{\ell'}^\top (\zeta_{\ell'} (\mathbf{X}_{\ell'}^\top \mathbf{w}_\ell + \mathbf{1}_{\ell'} b_\ell) + \boldsymbol{\xi}_{\ell'} - \mathbf{1}_{\ell'}) - \boldsymbol{\mu}_{\ell'}^\top \boldsymbol{\xi}_{\ell'} \end{aligned} \quad (9.36)$$

$\forall \ell \in \{+, -\}$, donde $\boldsymbol{\alpha}_{\ell'}$, $\boldsymbol{\mu}_{\ell'}$, $\mathbb{R}^{N'_\ell}$ son vectores con los valores de los multiplicadores de Lagrange. Su correspondiente conjunto de condiciones KKT son expresadas por:

$$\boldsymbol{\alpha}_{\ell'} \geq \mathbf{0}, \quad (9.37)$$

$$\zeta_{\ell'} (\mathbf{X}_{\ell'}^\top \mathbf{w}_\ell + \mathbf{1}_{\ell'} b_\ell) + \boldsymbol{\xi}_{\ell'} \geq \mathbf{1}_{\ell'}, \quad (9.38)$$

$$\boldsymbol{\alpha}_{\ell'}^\top (\zeta_{\ell'} (\mathbf{X}_{\ell'}^\top \mathbf{w}_\ell + \mathbf{1}_{\ell'} b_\ell) + \boldsymbol{\xi}_{\ell'} - \mathbf{1}_{\ell'}) = 0, \quad (9.39)$$

$$\boldsymbol{\mu}_{\ell'} \geq \mathbf{0}, \quad (9.40)$$

$$\boldsymbol{\xi}_{\ell'} \geq \mathbf{0}, \quad (9.41)$$

$$\boldsymbol{\mu}_{\ell'} \boldsymbol{\xi}_{\ell'} = \mathbf{0}, \quad (9.42)$$

$\forall \ell \in \{+, -\}$. Ahora hallamos las derivadas sobre \mathbf{w}_ℓ , b_ℓ y $\boldsymbol{\xi}_{\ell'}$, como:

$$\frac{\partial \mathcal{L}_\ell}{\partial \mathbf{w}_\ell} = \mathbf{0} \Rightarrow \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{w}_\ell + \mathbf{X}_\ell \mathbf{1}_\ell b_\ell + \zeta_{\ell'} \mathbf{X}_{\ell'} \boldsymbol{\alpha}_{\ell'} = \mathbf{0}, \quad (9.43)$$

$$\frac{\partial \mathcal{L}_\ell}{\partial b_\ell} = 0 \Rightarrow \mathbf{1}_\ell^\top \mathbf{X}_\ell^\top \mathbf{w}_\ell + \mathbf{1}_\ell^\top \mathbf{1}_\ell b_\ell + \zeta_{\ell'} \mathbf{1}_{\ell'}^\top \boldsymbol{\alpha}_{\ell'} = 0, \quad (9.44)$$

$$\frac{\partial \mathcal{L}_\ell}{\partial \boldsymbol{\xi}_{\ell'}} = 0 \Rightarrow c_{\ell',1} \mathbf{1}_{\ell'} - \boldsymbol{\alpha}_{\ell'} - \boldsymbol{\mu}_{\ell'} = \mathbf{0}, \quad (9.45)$$

Luego, combinando las ecuaciones (9.43) y (9.44) llegamos a:

$$\left(\begin{bmatrix} \mathbf{X}_\ell \\ \mathbf{1}_\ell^\top \end{bmatrix} \begin{bmatrix} \mathbf{X}_\ell^\top & \mathbf{1}_\ell \end{bmatrix} \right) \begin{bmatrix} \mathbf{w}_\ell \\ b_\ell \end{bmatrix} + \zeta_{\ell'} \begin{bmatrix} \mathbf{X}_{\ell'} \\ \mathbf{1}_{\ell'}^\top \end{bmatrix} \boldsymbol{\alpha}_{\ell'} = \mathbf{0}, \quad (9.46)$$

definiendo la matriz extendida $\mathbf{S}_\ell = \begin{bmatrix} \mathbf{X}_\ell^\top & \mathbf{1}_\ell \end{bmatrix}^\top$ y $\mathbf{z}_\ell = \begin{bmatrix} \mathbf{w}_\ell^\top & b_\ell \end{bmatrix}^\top$, y llevando la ecuación (9.46) en función de estas obtenemos que

$$\mathbf{S}_\ell^\top \mathbf{S}_\ell \mathbf{z}_\ell + \zeta_{\ell'} \mathbf{S}_{\ell'}^\top \boldsymbol{\alpha}_{\ell'} = \mathbf{0} \Rightarrow \mathbf{z}_{\ell'} = \zeta_{\ell'} \left(\mathbf{S}_\ell^\top \mathbf{S}_\ell \right)^{-1} \mathbf{S}_{\ell'}^\top \boldsymbol{\alpha}_{\ell'}. \quad (9.47)$$

$\forall \ell \in \{+, -\}$. Es importante notar que $\mathbf{S}_\ell^\top \mathbf{S}_\ell$ es siempre semi-definida positiva; sin embargo, es posible que no pueda estar bien condicionada en algunas ocasiones. Por lo anterior, se introduce un término de regularización del enfoque de regresión rígida tal como en [58], $\epsilon \mathbf{I}$ con $\epsilon \geq 0$, en aras de contrarrestar posibles problemas de mal condicionamiento de $\mathbf{S}_\ell^\top \mathbf{S}_\ell$. Aquí, \mathbf{I} es una matriz identidad de tamaño necesario. Entonces introduciendo este termino de regularización, la ecuación (9.47) toma la forma:

$$\mathbf{z}_{\ell'} = \zeta_{\ell'} \left(\mathbf{S}_\ell^\top \mathbf{S}_\ell + \epsilon \mathbf{I} \right)^{-1} \mathbf{S}_{\ell'}^\top \boldsymbol{\alpha}_{\ell'}, \quad (9.48)$$

Usando las definiciones de \mathbf{S}_ℓ y de \mathbf{z}_ℓ reescribimos el Lagrangiano en la ecuación (9.36) en función de estas dos variables. Luego, reemplazamos en esta la igualdad de la ecuación (9.47), si se requiere usar el termino de regularización ϵ por mal condicionamiento se reemplaza la ecuación (9.48), y obtenemos el problema de optimización dual:

$$\begin{aligned} \tilde{\boldsymbol{\alpha}}_{\ell'} &= \arg \max_{\boldsymbol{\alpha}_{\ell'}} \mathbf{1}_{\ell'}^\top \boldsymbol{\alpha}_{\ell'} - \frac{1}{2} \boldsymbol{\alpha}_{\ell'}^\top \mathbf{S}_{\ell'}^\top \left(\mathbf{S}_\ell \mathbf{S}_\ell^\top \right)^{-1} \mathbf{S}_{\ell'}^\top \boldsymbol{\alpha}_{\ell'} \\ \text{s.t. } &\mathbf{0} \leq \boldsymbol{\alpha}_{\ell'} \leq c_{\ell',1}. \end{aligned} \quad (9.49)$$

Como $\boldsymbol{\mu}_{\ell'}, \boldsymbol{\xi}_{\ell'} \geq 0$ y teniendo en consideración la ecuación (9.45), tenemos la restricción de los problemas duales expresados en la ecuación (9.49). Podemos observar que los problemas duales toman la forma de dos QPPs, más pequeñas que la QPP del entrenamiento de la SVM. Observemos que cada una de las QPPs de la TWSVM tiene una complejidad computacional de alrededor $\mathcal{O}(N_\ell^3) \forall \ell \in \{+, -\}$. Supongamos que $N_+ = N_- = N/2$, entonces tenemos que la complejidad computacional del TWSVM es de $\mathcal{O}(N^3/4)$ [57].

Ahora, una vez calculado los multiplicadores de Lagrange $\tilde{\boldsymbol{\alpha}}_{\ell'}$ usamos la ecuación (9.47) ó 9.48, dependiendo del acondicionamiento de las matriz $\mathbf{S}_\ell \mathbf{S}_\ell^\top$. Con lo anterior, calculamos el vector $\mathbf{z}_\ell = \begin{bmatrix} \tilde{\mathbf{w}}_\ell^\top & \tilde{b}_\ell \end{bmatrix}^\top$, y extraemos los valores para $\tilde{\mathbf{w}}_\ell$ y de $\tilde{b}_\ell \forall \ell \in \{+, -\}$. Es importante recalcar que para calcular, por ejemplo, los parámetros del hiperplano de la clase +1 es necesario encontrar los multiplicadores de Lagrange de la clase contraria, $\boldsymbol{\alpha}_-$. De lo anterior, podemos decir que los vectores de soporte de la clase ζ_ℓ provienen de las muestras de la clase $\zeta_{\ell'}$.

Por último, la etiqueta de una nueva muestra, $\mathbf{x}_* \in \mathbb{R}^P$, se asigna según a que hiperplano esta más cercana, haciendo uso de la distancia perpendicular. Lo anterior se resumen en la siguiente expresión como:

$$t_* = \arg \min_{\ell \in \{+, -\}} \frac{|f_\ell(\mathbf{x}_*)|}{\|\mathbf{w}_\ell\|_2}. \quad (9.50)$$

9.3.2. TWSVM No-lineal

La TWSVM extiende sus resultados a una extensión no-lineal, considerando la generación de dos superficies de separación en el espacio de entrada, construidas por alguna función kernel, en vez de hiperplanos. Las superficies de separación tienen la forma:

$$f_\ell(\mathbf{x}) = \mathbf{k}(\mathbf{x})^\top \mathbf{w}_\ell + b_\ell = 0, \quad \forall \ell \in \{+, -\}, \quad (9.51)$$

donde $\mathbf{k}(\cdot) : \mathbb{R}^P \rightarrow \mathbb{R}^N$ es un vector con elementos $k_n = \kappa(\mathbf{x}_n, \cdot) \quad \forall n \in \{1, 2, \dots, N\}$, además $\mathbf{x}_n \in \widehat{\mathbf{X}} = [\mathbf{X}_+ \ \mathbf{X}_-]$ y $\kappa(\cdot) : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$ es una función kernel escogida con anticipación. Nótese, que para la versión no-lineal del TWSVM $\mathbf{w}_+, \mathbf{w}_- \in \mathbb{R}^N$. Dicho lo anterior, los problemas de optimización primales toman la forma:

$$\begin{aligned} \tilde{\mathbf{w}}_+, \tilde{b}_+, \tilde{\xi}_- &= \arg \min_{\mathbf{w}_+, b_+, \xi_-} \frac{1}{2} \|\mathbf{K}_+^\top \mathbf{w}_+ + \mathbf{1}_+ b_+\|_2^2 + c_{-,1} \mathbf{1}_-^\top \xi_- \\ \text{s.t.} \quad &-1 \left(\mathbf{K}_-^\top \mathbf{w}_+ + \mathbf{1}_- b_+ \right) + \xi_- \geq \mathbf{1}_- \\ &\xi_- \geq 0, \end{aligned} \quad (9.52)$$

$$\begin{aligned} \tilde{\mathbf{w}}_-, \tilde{b}_-, \tilde{\xi}_+ &= \arg \min_{\mathbf{w}_-, b_-, \xi_+} \frac{1}{2} \|\mathbf{K}_-^\top \mathbf{w}_- + \mathbf{1}_- b_-\|_2^2 + c_{+,1} \mathbf{1}_+^\top \xi_+ \\ \text{s.t.} \quad &+1 \left(\mathbf{K}_+^\top \mathbf{w}_- + \mathbf{1}_+ b_- \right) + \xi_+ \geq \mathbf{1}_+ \\ &\xi_+ \geq 0, \end{aligned} \quad (9.53)$$

donde $\mathbf{K}_\ell \in \mathbb{R}^{N \times N_\ell}$ con elementos $k_{nn'} = \kappa(\mathbf{x}_n, \mathbf{x}_{n'})$, tal que $\mathbf{x}_n \in \widehat{\mathbf{X}}$ y $\mathbf{x}_{n'} \in \mathbf{X}_\ell$. Siguiendo el mismo procedimiento que en el apartado 9.3.1 pero cambiando las matrices $\mathbf{X}_\ell \quad \forall \ell \in \{+, -\}$ por las matrices $\mathbf{K}_\ell \quad \forall \ell \in \{+, -\}$. En consecuencia, las matrices extendidas $\mathbf{S}_\ell \in \mathbb{R}^{(N+1) \times N_\ell}$ son calculadas de la forma $\mathbf{S}_\ell = [\mathbf{K}_\ell^\top \ \mathbf{1}_\ell]^\top \quad \forall \ell \in \{+, -\}$. Entonces, los problemas duales toman la misma forma que 9.49. Así mismo, la función de decisión de la ecuación (9.50) se reescribe como:

$$t_* = \arg \min_{\ell \in \{+, -\}} \frac{|f_\ell(\mathbf{x}_*)|}{\|\mathbf{K}^{1/2} \mathbf{w}_\ell\|_2}, \quad (9.54)$$

donde $\mathbf{K} \in \mathbb{R}^{N \times N}$ es la matriz de Gram de las muestras de entrenamiento, con elementos $k_{nn'} = \kappa(\mathbf{x}_n, \mathbf{x}_{n'})$; $n, n' \in \{1, 2, \dots, N\}$ y $\mathbf{x}_n, \mathbf{x}_{n'} \in \widehat{\mathbf{X}}$.

9.4. Máquinas de Vectores de Soporte Gemelas Delimitadas (TBSVM)

Debido a que, la solución de los problemas duales del entrenamiento del TWSVM (ver ecuación (9.49)) dependen de que las matrices $\mathbf{S}_+ \mathbf{S}_+^\top$ y $\mathbf{S}_- \mathbf{S}_-^\top$ sean invertibles, lo cual implicaría que se debe asumir que sean no-singulares. Sin embargo, este requisito no puede ser satisfecho siempre. Jayadeva en [57] propone

regularizar estas matrices por medio del termino $\epsilon \mathbf{I}$, $\epsilon \geq 0$, con el fin de contrarrestar el mal acondicionamiento de estas matrices. Aún así, la teoría dual en TWSVM no es perfecta [59]. Por lo anterior, Shao et al en [59] proponen algunas mejoras a la formulación del TWSVM, al cual nombraron *Twin Bounded Support Vector Machine*–(TBSVM). El TBSVM al igual que el TWSVM, busca construir dos hiperplanos no-paralelos, con la diferencia de que en los problemas de optimización (ver ecuaciones (9.34) y (9.35) y ecuaciones (9.52) y (9.53)) , el TBSVM agrega un término de regularización con la idea de maximizar la margen de los hiperplanos, por medio de minimizar $\|\mathbf{w}_\ell\|_2^2$ y b_ℓ^2 para todo $\ell \in \{+, -\}$.

9.4.1. TBSVM lineal

Al igual que la versión lineal del TWSVM, el TBSVM busca construir dos hiperplanos no-paralelos, ver ecuación (9.33). Los problemas primales que propone TBSVM están dados por:

$$\begin{aligned} \tilde{\mathbf{w}}_\ell, \tilde{b}_\ell, \tilde{\boldsymbol{\xi}}_{\ell'} = \arg \min_{\mathbf{w}_+, b_+, \boldsymbol{\xi}_-} & \frac{c_{\ell,1}}{2} (\|\mathbf{w}_\ell\|_2^2 + b_\ell^2) + \frac{1}{2} \|\mathbf{X}_\ell^\top \mathbf{w}_\ell + \mathbf{1}_\ell b_\ell\|_2^2 + c_{\ell',2} \mathbf{1}_{\ell'}^\top \boldsymbol{\xi}_{\ell'} \\ \text{s.t.} & \quad \zeta_{\ell'} (\mathbf{X}_{\ell'}^\top \mathbf{w}_\ell + \mathbf{1}_{\ell'} b_\ell) + \boldsymbol{\xi}_{\ell'} \geq \mathbf{1}_{\ell'} \\ & \quad \boldsymbol{\xi}_{\ell'} \geq 0, \end{aligned} \quad (9.55)$$

$\forall \ell \in \{+, -\}$. Luego, el Lagrangiano para ambos problemas de optimización toma la forma:

$$\begin{aligned} \mathcal{L}_\ell(\mathbf{w}_\ell, b_\ell, \boldsymbol{\xi}_{\ell'}) = & \frac{c_{\ell,1}}{2} (\mathbf{q}_\ell^\top \mathbf{w}_\ell + b_\ell^2) + \frac{1}{2} (\mathbf{X}_\ell^\top \mathbf{w}_\ell + \mathbf{1}_\ell b_\ell)^\top (\mathbf{X}_\ell^\top \mathbf{w}_\ell + \mathbf{1}_\ell b_\ell) + c_{\ell',2} \mathbf{1}_{\ell'}^\top \boldsymbol{\xi}_{\ell'} - \dots \\ & \alpha_{\ell'}^\top (\zeta_{\ell'} (\mathbf{X}_{\ell'}^\top \mathbf{w}_\ell + \mathbf{1}_{\ell'} b_\ell) + \boldsymbol{\xi}_{\ell'} - \mathbf{1}_{\ell'}) - \boldsymbol{\mu}_{\ell'}^\top \boldsymbol{\xi}_{\ell'}, \end{aligned} \quad (9.56)$$

con $c_{\ell,1}, c_{\ell',2} \in \mathbb{R}^+$ como parámetros de regularización. El primer parámetro de regularización, penaliza el tamaño de los parámetros de los hiperplanos y el segundo regulariza las variables auxiliares, que indican cuando las muestras quedan mal clasificadas. El conjunto de condiciones KKT para los problemas de optimización del TBSVM son las mismos a las del TWSVM expresados en las ecuaciones (9.37) a (9.42). Ahora, hallamos las derivadas sobre \mathbf{w}_ℓ , b_ℓ y $\boldsymbol{\xi}_{\ell'}$, de lo cual obtenemos:

$$\frac{\partial \mathcal{L}_\ell}{\partial \mathbf{w}_\ell} = 0 \Rightarrow c_{\ell,1} \mathbf{w}_\ell + \mathbf{X}_\ell \mathbf{X}_\ell^\top \mathbf{w}_\ell + \mathbf{X}_\ell \mathbf{1}_\ell b_\ell + \zeta_{\ell'} \mathbf{X}_{\ell'} \boldsymbol{\alpha}_{\ell'} = 0, \quad (9.57)$$

$$\frac{\partial \mathcal{L}_\ell}{\partial b_\ell} = 0 \Rightarrow c_{\ell,1} b_\ell + \mathbf{1}_\ell^\top \mathbf{X}_\ell^\top \mathbf{w}_\ell + \mathbf{1}_\ell^\top \mathbf{1}_\ell b_\ell + \zeta_{\ell'} \mathbf{1}_{\ell'}^\top \boldsymbol{\alpha}_{\ell'} = 0, \quad (9.58)$$

$$\frac{\partial \mathcal{L}_\ell}{\partial \boldsymbol{\xi}_{\ell'}} = 0 \Rightarrow c_{\ell',2} \mathbf{1}_{\ell'} - \boldsymbol{\alpha}_{\ell'} - \boldsymbol{\mu}_{\ell'} = 0. \quad (9.59)$$

Siguiendo con el procedimiento, combinamos las ecuaciones (9.57) y (9.58) para obtener el siguiente sistema de ecuaciones:

$$\left(c_{\ell,1} \mathbf{I} + \begin{bmatrix} \mathbf{X}_\ell \\ \mathbf{1}_\ell^\top \end{bmatrix} \begin{bmatrix} \mathbf{X}_\ell^\top & \mathbf{1}_\ell \end{bmatrix} \right) \begin{bmatrix} \mathbf{w}_\ell \\ b_\ell \end{bmatrix} + \zeta_{\ell'} \begin{bmatrix} \mathbf{X}_{\ell'} \\ \mathbf{1}_{\ell'}^\top \end{bmatrix} \boldsymbol{\alpha}_{\ell'} = 0, \quad (9.60)$$

definimos la matriz extendida $\mathbf{S}_\ell = [\mathbf{X}_\ell^\top \quad \mathbf{1}_\ell]^\top$ y $\mathbf{z}_\ell = [\mathbf{w}_\ell^\top \quad b_\ell]^\top$, y reformulando la ecuación (9.60) en función de estas, y obtenemos:

$$\left(c_{\ell,1} \mathbf{I} + \mathbf{S}_\ell^\top \mathbf{S}_\ell \right) \mathbf{z}_\ell + \zeta_{\ell'} \mathbf{S}_{\ell'} \boldsymbol{\alpha}_{\ell'} = 0 \Rightarrow \mathbf{z}_{\ell'} = \zeta_{\ell'} \left(c_{1,\ell} \mathbf{I} + \mathbf{S}_\ell^\top \mathbf{S}_\ell \right)^{-1} \mathbf{S}_{\ell'} \boldsymbol{\alpha}_{\ell'}. \quad (9.61)$$

Usando las definiciones de S_ℓ y de z_ℓ reescribimos el Lagrangiano (de la ecuación (9.56)) en función de estas dos variables. Luego, reemplazamos en el Lagrangiano resultante la igualdad de la ecuación (9.61), y así se obtiene el problema de optimización dual de la forma:

$$\begin{aligned} \tilde{\alpha}_{\ell'} = \arg \max_{\alpha_{\ell'}} & \mathbf{1}_{\ell'}^\top \alpha_{\ell'} - \frac{1}{2} \alpha_{\ell'}^\top S_{\ell'}^\top (c_{\ell,1} \mathbf{I} + S_\ell S_\ell^\top)^{-1} S_{\ell'}^\top \alpha_{\ell'} \\ \text{s.t. } & \mathbf{0} \leq \alpha_{\ell'} \leq c_{\ell',2}. \end{aligned} \quad (9.62)$$

Igual que el TWSVM, para clasificar una nueva muestra x_* se tienen como criterio de decisión la mínima distancia entre el punto y los dos hiperplanos. Se usa la misma expresión de decisión que el TWSVM expresada en la ecuación (9.50).

9.4.2. TBSVM No-Lineal

De la misma forma que lo hizo las TWSVM, las TBSVM propuso una extensión no-lineal en donde en vez de construir dos hiperplanos, se generan dos superficies de separación generadas por una función kernel, y están expresadas por la ecuación 9.51.

Las funciones objetivos del TBSVM no-lineal son similares a las del TWSVM no-lineal, con la diferencia del termino de regularización. Dicho lo anterior, los problemas de optimización primales para el TBSVM no-lineal estan dados por:

$$\begin{aligned} \tilde{w}_\ell, \tilde{b}_\ell, \tilde{\xi}_{\ell'} = \arg \min_{w_\ell, b_\ell, \xi_{\ell'}} & \frac{c_{\ell,1}}{2} (\|w_\ell\|_2^2 + b_\ell^2) + \frac{1}{2} \|K_\ell^\top w_\ell + \mathbf{1}_\ell b_\ell\|_2^2 + c_{\ell',2} \mathbf{1}_{\ell'}^\top \xi_{\ell'} \\ \text{s.t. } & \xi_{\ell'} (K_{\ell'}^\top w_\ell + \mathbf{1}_{\ell'} b_\ell) + \xi_{\ell'} \geq \mathbf{1}_{\ell'} \\ & \xi_{\ell'} \geq 0, \end{aligned} \quad (9.63)$$

$\forall \ell \in \{+, -\}$. Siguiendo el mismo procedimiento que en el apartado 9.4.1 con las mismas definiciones usadas en el apartado 9.3.2. Además, reemplazamos las matrices X_ℓ por las matrices K_ℓ , $\forall \ell \in \{+, -\}$, previamente definidos en las secciones anteriores. Obtenemos las mismas formas de los problemas duales del TBSVM lineal, ver ecuación (9.62), con la excepción que la matrix extendida $S_\ell = [K_\ell^\top \quad \mathbf{1}_\ell]^\top$. Por último, la clasificación de una nueva muestra x_* se da por medio de la función de decisión 9.54.

9.5. Máquinas de vectores de soporte con un Lagrangiano ponderado

Como se ha expresado, las TWSVM y TBSVM tienen una complejidad computacional menor que las SVM estándar. Pero estos clasificadores poseen los mismos problemas ante bases de datos desbalanceadas. Como medida para contrarrestar estos problemas Shao et al. en [30] proponen *Weighted Lagrangian Twin Support Vector Machine* (WLTSVM), el cual al igual que las TWSVM busca construir dos hiperplanos no-paralelos. Además, ellos utilizan diferentes conjuntos de entrenamiento para generar ambos hiperplanos. El WLTSVM tiene como propiedades: (1) un submuestreo basado en grafos para reducir eficientemente el número de muestras de la clase de mayor muestras (clase mayoritaria), y el cual es robusto a muestras outliers; (2) la introducción analítica de un termino de sesgo ponderado, con el objetivo de mejorar el rendimiento sobre la clase de menor muestras (clase minoritaria), la cual tiende a ser ignorada por los clasificadores que no consideran el desbalance [6]; (3) los QPPs de este con respecto a los del TWSVM y TBSVM son diferentes, siendo más rápidos con rendimientos similares sobre base de datos sintéticas y del estado-del-arte.

9.5.1. WLTSVM lineal

Es importante recalcar que el WLTSVM fue diseñado para clasificación con bases de datos desbalanceadas, y por convención para estos problemas se toma la etiqueta +1 para la clase minoritaria (la de menor número de muestras) y la etiqueta -1 para la clase mayoritaria (la de mayor número de muestras). Dicho esto, la versión lineal del WLTSVM busca dos hiperplanos de la forma 9.33. Antes de definir el problema primal para el hiperplano minoritario (para la clase +1), debemos de definir el submuestreo basado en ν -nearest neighbor- $(\nu$ -NN), el cual busca caracterizar la densidad intra-clase de la clase mayoritaria (clase -1). En el grafo se considera que una muestra pertenece a la densidad intra-clase, cuando esta, junto a otra, son vecinos mutuos, ver figura 9.5(a). Esto implica que los puntos en las regiones de alta densidad tienen una mayor probabilidad de ser escogidos, mientras que los puntos que estén dentro de las regiones de baja densidad, por ejemplo, outliers, tendrán una baja probabilidad de ser seleccionados para el entrenamiento. Primero, definimos la matriz adyacente U , con elementos:

$$U_{ij} = \begin{cases} \rho, & \text{Si } \mathbf{x}_i \in N_\nu(j) \text{ y } \mathbf{x}_j \in N_\nu(i); \\ 0, & \text{En otro caso,} \end{cases} \quad (9.64)$$

donde $\mathbf{x}_i, \mathbf{x}_j \in \mathbf{X}_-$, $N_\nu(j)$ denota el conjunto de ν -vecinos más cercanos, dentro de la clase mayoritaria, de la muestra \mathbf{x}_j , y $\rho \in \mathbb{R}^+$ es un escalar escogido con anticipación, $i, j \in \{1, 2, \dots, N_-\}$. Entonces definimos el vector de coeficientes de submuestreo, \mathbf{u} con elementos

$$u_i = \begin{cases} 1, & \text{Si } \sum_j U_{ij} \geq \nu; \\ 0, & \text{En otro caso.} \end{cases} \quad (9.65)$$

Las muestras \mathbf{x}_i para la cual su correspondiente u_i es diferente de cero puede ser seleccionado, así obtenemos la matriz de muestras mayoritaria sub-muestreadas $\mathbf{X}'_- \in \mathbb{R}^{P \times N'_-}$, donde N'_- es el número de muestras de la clase mayoritaria escogidas por el submuestreo y $N_- \geq N'_-$. Fijémonos, que no hay forma de controlar el nivel de submuestreo, por esta razón es posible que no se elimine el desbalance entre ambas clases y hasta que llegue a suceder que $N_+ > N'_-$. De forma consecutiva, los autores introducen unos pesos en la construcción del hiperplano minoritario. Muy similar al WSVM, se define la matriz de pesos diagonal $\mathbf{D}_- = \min\{N_+/N'_-, N'_-/N_+\} \mathbf{I}$, con $\mathbf{D}_- \in \mathbb{R}^{N'_- \times N'_-}$.

Ahora, usando la matriz de muestras sub-muestreadas \mathbf{X}'_- y la diagonal de pesos \mathbf{D}_- , encontraríamos los parámetros del hiperplano minoritario solucionando el problema de optimización:

$$\begin{aligned} \tilde{\mathbf{w}}_+, \tilde{b}_+, \tilde{\boldsymbol{\xi}}_- &= \arg \min_{\mathbf{w}_+, b_+, \boldsymbol{\xi}_-} \frac{c_{+,1}}{2} \left(\|\mathbf{w}_+\|_2^2 + b_+^2 + \boldsymbol{\xi}_-^\top \mathbf{D}_- \boldsymbol{\xi}_- \right) + \frac{1}{2} \|\mathbf{X}'_+^\top \mathbf{w}_+ + \mathbf{1}_+ b_+\|_2^2 \\ \text{s.t.} \quad & -1 \left(\mathbf{X}'_-^\top \mathbf{w}_+ + \mathbf{1}_- b_+ \right) + \boldsymbol{\xi}_- \geq \mathbf{1}_- \\ & \boldsymbol{\xi}_- \geq 0, \end{aligned} \quad (9.66)$$

donde $\boldsymbol{\xi}_- \in \mathbb{R}^{N'_-}$ son las variables auxiliares de las muestras de la clase mayoritaria sub-muestreadas y $c_{+,1} \in \mathbb{R}^+$.

Para el hiperplano mayoritario, se construye un grafo que caracterice la separabilidad inter-clase entre las muestras de la clase mayoritaria de la minoritaria. En el grafo inter-clase, un par de muestras son conectadas si ellas provienen de la clase minoritaria y mayoritaria, respectivamente, y la muestra de la clase mayoritaria es un ν' -NN de la muestra minoritaria (ver figura 9.5(b)). Esto implica que las muestras en la región marginal tienen más posibilidades de ser escogidos. Específicamente, la matriz adyacente al hiperplano mayoritario es denotada por $\mathbf{V} \in \mathbb{R}^{N_- \times N_+}$, con elementos:

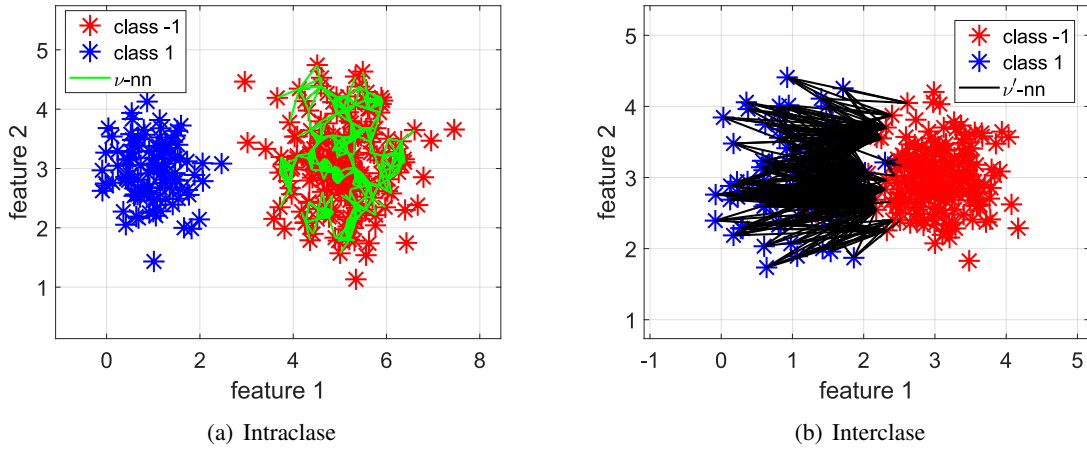


Figura 9.5.: a) Densidad intraclase según el submuestreo basado en ν -nn. b) Densidad interclase según el submuestreo basado en ν' -nn. Para la generación de ambas gráficas fijamos $\nu=\nu'=6$.

$$V_{ij} = \begin{cases} \rho', & \text{Si } \mathbf{x}_i \in N_{\nu'}(j); \\ 0, & \text{En otro caso,} \end{cases} \quad (9.67)$$

donde $N_{\nu'}(j)$ denota el conjunto de muestras compuestas por los ν' -vecinos más cercanos, pertenecientes a la clase mayoritaria, de la muestra $\mathbf{x}_j \in \mathbf{X}_+$, y $\rho' \in \mathbb{R}^+$ es un escalar, $i \in \{1, 2, \dots, N_-\}$ y $j \in \{1, 2, \dots, N_+\}$. Entonces, se define el vector de coeficientes de submuestreo $\mathbf{v} \in \mathbb{R}^{N_-}$ con elementos:

$$v_i = \begin{cases} 1, & \text{Si } \sum_j V_{ij} \geq \nu'; \\ 0, & \text{En otro caso.} \end{cases} \quad (9.68)$$

Luego construimos la matriz sub-muestreada para el hiperplano mayoritario $\bar{\mathbf{X}}_- \in \mathbb{R}^{P \times N'_-}$, $\bar{\mathbf{X}}_- = [\mathbf{x}_i \mid v_i = 1]$, con $N'_- \leq N_-$ y $i \in 1, 2, \dots, N_-$. De forma similar se define la matriz diagonal de pesos $\mathbf{D}_+ \in \mathbb{R}^{N_+ \times N_+}$, dada por $\mathbf{D}_+ = \min\{N_+/N'_-, N'_-/N_+\} \mathbf{I}$. Entonces, los parámetros del hiperplano mayoritarios se encontrarían resolviendo el problema de optimización:

$$\begin{aligned} \tilde{\mathbf{w}}_-, \tilde{b}_-, \tilde{\boldsymbol{\xi}}_+ &= \arg \min_{\mathbf{w}_-, b_-, \boldsymbol{\xi}_+} \frac{c_{+,1}}{2} (\|\mathbf{w}_-\|_2^2 + b_-^2 + \boldsymbol{\xi}_+^\top \mathbf{D}_+ \boldsymbol{\xi}_+) + \frac{1}{2} \|\bar{\mathbf{X}}_-^\top \mathbf{w}_- + \mathbf{1}_- b_- \|_2^2 \\ \text{s.t. } & (\mathbf{X}_+^\top \mathbf{w}_- + \mathbf{1}_+ b_-) + \boldsymbol{\xi}_+ \geq \mathbf{1}_+ \\ & \boldsymbol{\xi}_+ \geq \mathbf{0}, \end{aligned} \quad (9.69)$$

Para resolver los problemas de optimización expresados en 9.66 y 9.69, son llevados a una forma dual. Primero, el Lagrangiano de 9.66 esta dado por

$$\begin{aligned} \mathcal{L}_+(\mathbf{w}_+, b_+, \boldsymbol{\xi}_-, \boldsymbol{\alpha}_-) &= \frac{c_{+,1}}{2} (\mathbf{w}_+^\top \mathbf{w}_+ + b_+^2) + \frac{1}{2} (\mathbf{X}_+^\top \mathbf{w}_+ + \mathbf{1}_+ b_+)^\top (\mathbf{X}_+^\top \mathbf{w}_+ + \mathbf{1}_+ b_+) + \dots \\ &\quad \dots + c_{+,1} \boldsymbol{\xi}_-^\top \mathbf{D}_- \boldsymbol{\xi}_- - \boldsymbol{\alpha}_-^\top (-1 (\mathbf{X}_-^\top \mathbf{w}_+ + \mathbf{1}_- b_+) + \boldsymbol{\xi}_- - \mathbf{1}_-) \end{aligned} \quad (9.70)$$

donde $\boldsymbol{\alpha}_- \in \mathbb{R}^{N'_-}$ es el vector con los multiplicadores de Lagrange de las muestras sub-muestreadas para \mathbf{X}'_- . Las condiciones de KKT del problema 9.66 están dadas por:

$$-1(\mathbf{X}'_{-}\mathbf{w}_{+} + \mathbf{1}_{+}b_{+}) + \xi_{-} - \mathbf{1}_{-} \geq 0 \quad (9.71)$$

$$\alpha_{-}^{\top}(-1(\mathbf{X}'_{-}\mathbf{w}_{+} + \mathbf{1}_{+}b_{+}) + \xi_{-} - \mathbf{1}_{-}) = 0 \quad (9.72)$$

$$\alpha_{-} \geq 0 \quad (9.73)$$

$$\xi_{-} \geq 0. \quad (9.74)$$

Las derivadas con respecto a \mathbf{w}_{+} , b_{+} y ξ_{-} toman la forma:

$$\frac{\partial \mathcal{L}_{+}}{\partial \mathbf{w}_{+}} = 0 \Rightarrow c_{+,1}\mathbf{w}_{+} + \mathbf{X}_{+}\mathbf{X}_{+}^{\top}\mathbf{w}_{+} + \mathbf{X}_{+}\mathbf{1}_{+}b_{+} + \mathbf{X}'_{-}\alpha_{-} = 0 \quad (9.75)$$

$$\frac{\partial \mathcal{L}_{+}}{\partial b_{+}} = 0 \Rightarrow c_{+,1}b_{+} + \mathbf{1}_{+}^{\top}\mathbf{X}_{+}^{\top}\mathbf{w}_{+} + \mathbf{1}_{+}^{\top}\mathbf{1}_{+}b_{+} + \mathbf{1}_{-}^{\top}\alpha_{-} = 0 \quad (9.76)$$

$$\frac{\partial \mathcal{L}_{+}}{\partial \xi_{-}} = 0 \Rightarrow c_{+,1}\mathbf{D}_{-}\xi_{-} - \alpha_{-} = 0. \quad (9.77)$$

Ahora, combinando las ecuaciones (9.75) y (9.76) nos conduce al siguiente sistema de ecuaciones, dado por:

$$\left(c_{+,1}\mathbf{I} + \begin{bmatrix} \mathbf{X}_{+} \\ \mathbf{1}_{+}^{\top} \end{bmatrix} \begin{bmatrix} \mathbf{X}_{+}^{\top} & \mathbf{1}_{+} \end{bmatrix} \right) \begin{bmatrix} \mathbf{w}_{+} \\ b_{+} \end{bmatrix} + \begin{bmatrix} \mathbf{X}'_{-} \\ \mathbf{1}_{-}^{\top} \end{bmatrix} \alpha_{-} = 0 \quad (9.78)$$

usando la definición de matriz extendida \mathbf{S}_{+} usada en el apartado 9.3, definimos la matriz extendida $\mathbf{S}'_{-} = [\mathbf{X}'_{-}^{\top} \quad \mathbf{1}_{-}^{\top}]^{\top}$ y $\mathbf{z}_{+} = [\mathbf{w}_{+}^{\top} \quad b_{+}]^{\top}$. Usamos estas definiciones sobre la ecuación (9.78) y despejando \mathbf{z}_{+} del resultado, haremos que:

$$\mathbf{z}_{+} = -(\mathbf{S}_{+}\mathbf{S}_{+}^{\top} + c_{+,1}\mathbf{I})^{-1}\mathbf{S}'_{-}\alpha_{-}. \quad (9.79)$$

Usando las definiciones antes mencionadas sobre el Lagrangeano 9.70 y sustituyendo la igualdad 9.79 en la misma, y usando la ecuación (9.77) obtenemos el problema dual de 9.66 como sigue:

$$\begin{aligned} \tilde{\alpha}_{-} = \arg \max_{\alpha_{-}} & -\frac{1}{2}\alpha_{-}^{\top} \left(\mathbf{S}'_{-}^{\top} (\mathbf{S}_{+}\mathbf{S}_{+}^{\top} + c_{+,1}\mathbf{I})^{-1} \mathbf{S}'_{-} + \frac{1}{c_{+,1}}\mathbf{D}_{-}^{-1} \right) \alpha_{-} + \mathbf{1}_{-}^{\top} \alpha_{-} \\ \text{s.t. } & \alpha_{-} \geq 0. \end{aligned} \quad (9.80)$$

Siguiendo el mismo procedimiento para obtener el problema dual del hiperplano de la clase minoritaria, resolvemos el problema de optimización 9.69 para entrenar el hiperplano de la clase mayoritaria, obteniendo:

$$\begin{aligned} \tilde{\alpha}_{+} = \arg \max_{\alpha_{+}} & -\frac{1}{2}\alpha_{+}^{\top} \left(\mathbf{S}_{+}^{\top} (\bar{\mathbf{S}}_{-}\bar{\mathbf{S}}_{-}^{\top} + c_{-,1}\mathbf{I})^{-1} \mathbf{S}_{+} + \frac{1}{c_{-,1}}\mathbf{D}_{+}^{-1} \right) \alpha_{+} + \mathbf{1}_{+}^{\top} \alpha_{+} \\ \text{s.t. } & \alpha_{+} \geq 0. \end{aligned} \quad (9.81)$$

donde $\alpha_{+} \in \mathbb{R}^{N_{+}}$ es el vector con los multiplicadores de Lagrange de las muestras de la clase +1 (minoritarias), $\bar{\mathbf{S}}_{-} = [\bar{\mathbf{X}}_{-}^{\top} \quad \mathbf{1}_{-}^{\top}]^{\top}$ y $\mathbf{z}_{-} = [\mathbf{w}_{-}^{\top} \quad b_{-}]^{\top}$; a partir de las condiciones de KKT obtenemos que:

$$\mathbf{z}_{-} = (\bar{\mathbf{S}}_{-}\bar{\mathbf{S}}_{-}^{\top} + c_{-,1}\mathbf{I})^{-1} \mathbf{S}_{+}\alpha_{+}. \quad (9.82)$$

Una vez se haya solucionado \mathbf{w}_{+} , b_{+} , \mathbf{w}_{-} y b_{-} , una nueva muestra $\mathbf{x}_{*} \in \mathbb{R}^P$ es clasificada por medio de la función de decisión 9.50.

9.5.2. WLTSVM No-Lineal

Al igual que el TWSVM y el TBSVM, el WLTSVM realiza una extensión para clasificación no-lineal, a través de la creación de superficies de separación generadas por una función kernel escogida con antelación. Estas superficies $f_+(x)=0$ y $f_-(x)=0$ están definidas en la ecuación (9.51).

Siguiendo con el mismo procedimiento que en su versión lineal y habiendo generado las matrices sub-muestreadas \mathbf{X}'_+ y \mathbf{X}'_- . Y definiendo la matriz kernel $\mathbf{K}'_- \in \mathbb{R}^{N \times N'_-}$ con elementos $k_{nn'} = \kappa(\mathbf{x}_n, \mathbf{x}_{n'})$ donde $\mathbf{x}_n \in \widehat{\mathbf{X}}$ y $\mathbf{x}_{n'} \in \mathbf{X}'_-$; y $\bar{\mathbf{K}}_- \in \mathbb{R}^{N \times N''_-}$ con elementos $k_{nn'} = \kappa(\mathbf{x}_n, \mathbf{x}_{n'})$ donde $\mathbf{x}_n \in \widehat{\mathbf{X}}$ y $\mathbf{x}_{n'} \in \bar{\mathbf{X}}_-$.

Los problemas de optimización toman la forma:

$$\begin{aligned} \tilde{w}_+, \tilde{b}_+, \tilde{\xi}_- &= \arg \min_{w_+, b_+, \xi_-} \frac{c_+, 1}{2} \left(\|w_+\|_2^2 + b_+^2 + \xi_-^\top \mathbf{D}_- \xi_- \right) + \frac{1}{2} \|\mathbf{K}'_+^\top w_+ + \mathbf{1}_+ b_+\|_2^2 \\ \text{s.t.} \quad &-1 \left(\mathbf{K}'_-^\top w_+ + \mathbf{1}_- b_+ \right) + \xi_- \geq \mathbf{1}_- \\ &\xi_- \geq 0, \end{aligned} \tag{9.83}$$

$$\begin{aligned} \tilde{w}_-, \tilde{b}_-, \tilde{\xi}_+ &= \arg \min_{w_-, b_-, \xi_+} \frac{c_-, 1}{2} \left(\|w_-\|_2^2 + b_-^2 + \xi_+^\top \mathbf{D}_+ \xi_+ \right) + \frac{1}{2} \|\bar{\mathbf{K}}_-^\top w_- + \mathbf{1}_- b_-\|_2^2 \\ \text{s.t.} \quad &\left(\mathbf{K}_+^\top w_- + \mathbf{1}_+ b_- \right) + \xi_+ \geq \mathbf{1}_+ \\ &\xi_+ \geq 0, \end{aligned} \tag{9.84}$$

Siguiendo el mismo procedimiento que en el apartado 9.5.1 podemos llegar a las mismas formas duales que en 9.80 y 9.81, con las matrices extendidas $\mathbf{S}_+ = [\mathbf{K}_+^\top \quad \mathbf{1}_+]^\top$, $\mathbf{S}'_- = [\mathbf{K}'_-^\top \quad \mathbf{1}_-]^\top$ y $\bar{\mathbf{S}}_- = [\bar{\mathbf{K}}_-^\top \quad \mathbf{1}_-]^\top$. También los vectores z_+ y z_- toman la misma forma que las ecuaciones (9.79) y (9.82), respectivamente. La clasificación de una nueva muestra \mathbf{x}_* esta dada por la función de decisión 9.54.

Parte II.

Propuesta

10. Máquinas de Vectores de Soporte Gemelas mejorada (ETWSVM)

Inducimos la solución de dos problemas claves que surgen cuando se tratan con datos desbalanceados dentro de modelos basados en TWSVM: *i)* la extensión no-lineal de la función de decisión a través de un mapeo adecuado a *espacios de Hilbert con kernel reproductivo*–(RKHS) con el fin de mejorar la representación de los datos de entrada; y *ii)* la estimación de los parámetros de la función kernel, por ejemplo, la matriz de covarianza de un kernel Gaussiano, para aprovechar la separabilidad de los datos y para evitar un conocimiento previo del usuario con respecto al ajuste de la función kernel.

Presentamos *máquinas de vectores de soporte gemelas mejorada*–(ETWSVM) para la clasificación binaria de datos desbalanceadas. Nuestro clasificador comprende una reforma a la formulación dual del problema de optimización cuadrática del TBSVM basándose del bien conocido lema de la matriz inversa, lo cual permite codificar una función de decisión en un RKHS con una dimensión posiblemente infinita. Además, una técnica basada en alineamiento de kernel es incluido durante el entrenamiento del ETWSVM como una estimación basada en datos de la función de similitud, kernel. Por último, nuestra propuesta es extendida a problemas multi-clase a través de las estrategias *uno-versus-Resto*–(OvR) y *uno-versus-uno*–(OvO) [39].

10.1. Fundamentos en ETWSVM

Nuestra propuesta básicamente es una reformulación de la versión no-lineal del TBSVM. Comenzando que las versiones no-lineales del TWSVM y del TBSVM buscan dos superficies de separación. Sugerimos que en vez de estas dos superficies se busque la construcción de dos hiperplanos no-paralelos en un espacio de RKHS de la formas

$$f_+(x) = \varphi(x)^\top w_+ + b_+ \quad \text{y} \quad f_-(x) = \varphi(x)^\top w_- + b_- \quad (10.1)$$

donde $w_\ell \in \mathcal{H}$ es el vector normal al hiperplano correspondiente, $b_\ell \in \mathbb{R}$ es el termino de sesgo del hiperplano correspondiente, para todo $\ell \in \{+, -\}$, y $\mathcal{H} \subset \mathbb{R}^Q$ es un RKHS. En general, \mathcal{H} se puede pensar como un espacio de características de alta-dimensionalidad (de “infinita” dimensionalidad, $Q \rightarrow \infty$) [11]. Al igual que el TBSVM versión lineal, apartado 9.4, se busca que cada hiperplano este lo más cerca posible de las muestras de su respectiva clase y lo más lejos posible de la otra clase y a la vez regularizar los parámetros de los hiperplanos. Dicho lo anterior, los problemas de optimización primal están dados por:

$$\begin{aligned} \tilde{w}_\ell, \tilde{b}_\ell, \tilde{\xi}_{\ell'} &= \arg \min_{w_+, b_+, \xi_-} \frac{c_{\ell,1}}{2} \left(\|w_\ell\|_2^2 + b_\ell^2 \right) + \frac{1}{2} \|\Phi_\ell^\top w_\ell + \mathbf{1}_\ell b_\ell\|_2^2 + c_{\ell',2} \mathbf{1}_{\ell'}^\top \xi_{\ell'} \\ \text{s.t.} \quad & \zeta_{\ell'} \left(\Phi_{\ell'}^\top w_\ell + \mathbf{1}_{\ell'} b_\ell \right) + \xi_{\ell'} \geq \mathbf{1}_{\ell'} \\ & \xi_{\ell'} \geq 0, \end{aligned} \quad (10.2)$$

para $\ell \in \{+, -\}$; donde $\Phi_\ell \in \mathbb{R}^{Q \times N_\ell}$ es la matriz de entrada con las muestras de la clase ζ_ℓ mapeadas por medio de la función de mapeo $\varphi(\cdot) : \mathbb{R}^P \rightarrow \mathcal{H}$. Es pertinente recordar que $\zeta_\ell = +1$ si $\ell = +$, y $\zeta_\ell = -1$ para el caso contrario. Dicho lo anterior, podemos definir el Lagrangiano para calcular el hiperplano $f_\ell(\mathbf{x})=0$ como:

$$\begin{aligned} \mathcal{L}_\ell(\mathbf{w}_\ell, b_\ell, \boldsymbol{\xi}_{\ell'}, \boldsymbol{\alpha}_{\ell'}, \boldsymbol{\mu}_{\ell'}) = & \frac{c_{\ell,1}}{2} (\mathbf{w}_\ell^\top \mathbf{w}_\ell + b_\ell^2) + \frac{1}{2} (\Phi_\ell^\top \mathbf{w}_\ell + \mathbf{1}_\ell b_\ell)^\top (\Phi_\ell^\top \mathbf{w}_\ell + \mathbf{1}_\ell b_\ell) + \dots \\ & \dots + c_{2,\ell'} \mathbf{1}_{\ell'}^\top \boldsymbol{\xi}_{\ell'} - \boldsymbol{\alpha}_{\ell'}^\top (\zeta_{\ell'} (\Phi_{\ell'}^\top \mathbf{w}_\ell + \mathbf{1}_{\ell'} b_\ell) + \boldsymbol{\xi}_{\ell'} - \mathbf{1}_{\ell'}) - \boldsymbol{\mu}_{\ell'}^\top \boldsymbol{\xi}_{\ell'}, \end{aligned} \quad (10.3)$$

donde $\boldsymbol{\xi}_{\ell'}, \boldsymbol{\mu}_{\ell'} \in \mathbb{R}^{N_{\ell'}}$ son los vectores con los multiplicadores de Lagrange de las condiciones expresadas en el problema de optimización 10.2 para las muestras de la clase $\zeta_{\ell'}$. Las condiciones KKT del problema de optimización están dados por:

$$\boldsymbol{\alpha}_{\ell'} \geq \mathbf{0}, \quad (10.4)$$

$$\zeta_{\ell'} (\Phi_{\ell'}^\top \mathbf{w}_\ell + \mathbf{1}_{\ell'} b_\ell) + \boldsymbol{\xi}_{\ell'} \geq \mathbf{1}_{\ell'}, \quad (10.5)$$

$$\boldsymbol{\alpha}_{\ell'}^\top (\zeta_{\ell'} (\Phi_{\ell'}^\top \mathbf{w}_\ell + \mathbf{1}_{\ell'} b_\ell) + \boldsymbol{\xi}_{\ell'} - \mathbf{1}_{\ell'}) = 0, \quad (10.6)$$

$$\boldsymbol{\mu}_{\ell'} \geq \mathbf{0}, \quad (10.7)$$

$$\boldsymbol{\xi}_{\ell'} \geq \mathbf{0}, \quad (10.8)$$

$$\boldsymbol{\mu}_{\ell'}^\top \boldsymbol{\xi}_{\ell'} = 0. \quad (10.9)$$

Luego, las derivadas con respecto a \mathbf{w}_ℓ , b_ℓ y $\boldsymbol{\xi}_{\ell'}$ toman la forma:

$$\frac{\partial \mathcal{L}_\ell}{\partial \mathbf{w}_\ell} = \mathbf{0} \Rightarrow c_{\ell,1} \mathbf{w}_\ell + \Phi_\ell \Phi_\ell^\top \mathbf{w}_\ell + \Phi_\ell \mathbf{1}_\ell b_\ell + \zeta_{\ell'} \Phi_{\ell'} \boldsymbol{\alpha}_{\ell'} = \mathbf{0}, \quad (10.10)$$

$$\frac{\partial \mathcal{L}_\ell}{\partial b_\ell} = 0 \Rightarrow c_{\ell,1} b_\ell + \mathbf{1}_\ell^\top \Phi_\ell^\top \mathbf{w}_\ell + \mathbf{1}_\ell^\top \mathbf{1}_\ell b_\ell + \zeta_{\ell'} \mathbf{1}_{\ell'}^\top \boldsymbol{\alpha}_{\ell'} = 0, \quad (10.11)$$

$$\frac{\partial \mathcal{L}_\ell}{\partial \boldsymbol{\xi}_{\ell'}} = 0 \Rightarrow c_{\ell',2} \mathbf{1}_{\ell'} - \boldsymbol{\alpha}_{\ell'} - \boldsymbol{\mu}_{\ell'} = \mathbf{0}. \quad (10.12)$$

Combinando las ecuaciones (10.10) y (10.11) con el fin de generar el sistema de ecuaciones,

$$\left(c_{\ell,1} \mathbf{I} + \begin{bmatrix} \Phi_\ell \\ \mathbf{1}_\ell^\top \end{bmatrix} \begin{bmatrix} \Phi_\ell^\top & \mathbf{1}_\ell \end{bmatrix} \right) \begin{bmatrix} \mathbf{w}_\ell \\ b_\ell \end{bmatrix} + \zeta_{\ell'} \begin{bmatrix} \Phi_{\ell'} \\ \mathbf{1}_{\ell'}^\top \end{bmatrix} \boldsymbol{\alpha}_{\ell'} = \mathbf{0}, \quad (10.13)$$

y definiendo las matrices extendidas $\mathbf{S}_\ell = \begin{bmatrix} \Phi_\ell^\top & \mathbf{1}_\ell \end{bmatrix}^\top$ y $\mathbf{z}_\ell = \begin{bmatrix} \mathbf{w}_\ell^\top & b_\ell \end{bmatrix}^\top$, para $\ell \in \{+, -\}$. Haciendo uso de las matrices definidas anteriormente, el sistema de ecuaciones 10.13 se reescribe de la siguiente forma:

$$\left(c_{\ell,1} \mathbf{I} + \mathbf{S}_\ell \mathbf{S}_\ell^\top \right) \mathbf{z}_\ell - \zeta_{\ell'} \mathbf{S}_{\ell'} \boldsymbol{\alpha}_{\ell'} = \mathbf{0}. \quad (10.14)$$

Ahora, despejamos \mathbf{z}_ℓ de la ecuación (10.14), y obtenemos que:

$$\mathbf{z}_\ell = \zeta_{\ell'} \left(c_{\ell,1} \mathbf{I} + \mathbf{S}_\ell \mathbf{S}_\ell^\top \right)^{-1} \mathbf{S}_{\ell'} \boldsymbol{\alpha}_{\ell'}. \quad (10.15)$$

Notemos que el calculo de la matrices de covarianza $\Phi_\ell \Phi_\ell^\top$, para $\ell \in \{+, -\}$, están en el espacio RKHS, las cuales no pueden ser calculadas debido a la “infinita”-dimensionalidad, propiedad de la función de mapeo $\varphi(\cdot)$. Sin embargo, aplicamos el lema de la matriz inversa, dado por:

$$(A + BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}, \quad (10.16)$$

y haciendo uso de las identificaciones:

$$c_{\ell,1}I \rightarrow A \quad S_{\ell} \rightarrow B \quad I \rightarrow C \quad S_{\ell}^{\top} \rightarrow D \quad (10.17)$$

obtenemos que:

$$(c_{\ell,1}I + S_{\ell}S_{\ell}^{\top})^{-1} = \frac{1}{2c_{\ell,1}}(I - S_{\ell}(c_{\ell,1}I + S_{\ell}S_{\ell})^{-1}S_{\ell}^{\top}S_{\ell})\alpha_{\ell'}, \quad (10.18)$$

para $\ell \in \{+, -\}$. Dicho lo anterior, llevamos el Lagrangiano (ecuación (10.3)) en función de las matrices S_{ℓ} y del vector z_{ℓ} , para $\ell \in \{+, -\}$, y reemplazando este último (z_{ℓ}) acorde a la ecuación (10.15), lo cual nos lleva a que el problema de optimización dual para el ETWSVM toma la forma:

$$\begin{aligned} \hat{\alpha}_{\ell'} = \arg \max_{\alpha_{\ell'}} & \mathbf{1}_{\ell'}^{\top} \alpha_{\ell'} - \frac{1}{2c_{\ell,1}} \alpha_{\ell'}^{\top} (S_{\ell'}^{\top} S_{\ell}' - S_{\ell'}^{\top} S_{\ell} (c_{\ell,1}I + S_{\ell}^{\top} S_{\ell})^{-1} S_{\ell}^{\top} S_{\ell}') \alpha_{\ell'} \\ \text{s.t. } & \mathbf{0} \leq \alpha_{\ell'} \leq c_{\ell',2} \mathbf{1}_{\ell'}. \end{aligned} \quad (10.19)$$

Basándonos en el producto interno en RKHS, las matrices $K_{\ell,\ell'}$, $\hat{K}_{\ell,\ell'} \in \mathbb{R}^{N_{\ell} \times N_{\ell'}}$ pueden ser inferidas como: $K_{\ell,\ell'} = \Phi_{\ell}^{\top} \Phi_{\ell'}$ y $\hat{K}_{\ell,\ell'} = S_{\ell}^{\top} S_{\ell'} = \hat{K}_{\ell,\ell'} + \mathbf{1}_{\ell} \mathbf{1}_{\ell'}^{\top}$, respectivamente. Es importante precisar que la matriz $K_{\ell,\ell'}$ tiene elementos $k_{nn'} = \kappa(\mathbf{x}_n, \mathbf{x}_{n'}) = \langle \varphi(\mathbf{x}_n), \varphi(\mathbf{x}_{n'}) \rangle_{\mathcal{H}}$, con $n \in \{1, \dots, N_{\ell}\}$ y $n' \in \{1, \dots, N_{\ell'}\}$, haciendo uso de lo expresado en el apartado 7.2, donde $\kappa : \mathbb{R}^P \times \mathbb{R}^P \rightarrow \mathbb{R}$ es la función kernel escogida con anticipación. En este sentido, el problema de optimización dual de la ecuación (10.19) queda de la forma

$$\begin{aligned} \hat{\alpha}_{\ell'} = \arg \max_{\alpha_{\ell'}} & \mathbf{1}_{\ell'}^{\top} \alpha_{\ell'} - \frac{1}{2c_{\ell',2}} \alpha_{\ell'}^{\top} (\hat{K}_{\ell',\ell'} - \hat{K}_{\ell',\ell} A_{\ell}^{-1} \hat{K}_{\ell,\ell'}) \alpha_{\ell'} \\ \text{s.t. } & \mathbf{0} \leq \alpha_{\ell'} \leq c_{\ell',2} \mathbf{1}_{\ell'}, \end{aligned} \quad (10.20)$$

donde $A_{\ell} = \hat{K}_{\ell,\ell} + c_{\ell,1}I$. La solución de los problemas de programación cuadrática en 10.20 (uno por cada $\alpha_{\ell'}, \forall \ell' \in \{+, -\}$), permiten obtener los vectores de pesos $z_{\ell} \in \mathbb{R}^{Q+1}$ en \mathcal{H} , El cual puede ser expresado como en la ecuación (10.15). De nuevo, el calculo de la matriz de covarianza en RKHS es eludido por medio del lema de la matriz inversa, tomando la forma:

$$z_{\ell} = \frac{\zeta_{\ell'}}{c_{\ell,1}} (S_{\ell'} - S_{\ell'} A_{\ell}^{-1} \hat{K}_{\ell,\ell'}) \alpha_{\ell'} \quad (10.21)$$

El puntaje de una nueva muestra, $\mathbf{x}_* \in \mathbb{R}^P$, para el hiperplano de la clase ζ_{ℓ} es calculado como $f_{\ell}(\mathbf{x}_*) = \varphi(\mathbf{x}_*)^{\top} \mathbf{w}_{\ell} + b_{\ell} = \mathbf{s}_*^{\top} z_{\ell}$, donde $\mathbf{s}_* = [\varphi(\mathbf{x}_*)^{\top} \quad 1]^{\top}$. Luego, aplicando la función kernel, este se reescribe de la forma:

$$f_{\ell}(\mathbf{x}_*) = \frac{\zeta_{\ell'}}{c_{\ell,1}} (\hat{\mathbf{k}}_{*,\ell'}^{\top} - \hat{\mathbf{k}}_{*,\ell}^{\top} A_{\ell}^{-1} \hat{K}_{\ell,\ell'}) \alpha_{\ell'}, \quad (10.22)$$

donde $\mathbf{k}_{*,\ell} \in \mathbb{R}^{N_{\ell}}$ contienen los elementos $\mathbf{k}_{*,\ell} = \{\kappa(\mathbf{x}_*, \mathbf{x}_n) : t_n = \zeta_{\ell}\}$, y $\hat{\mathbf{k}}_{*,\ell} = \mathbf{k}_{*,\ell} + \mathbf{1}_{\ell}$. Para encontrar la función de decisión, primero reescribimos z_{ℓ} como:

$$z_{\ell} = \begin{bmatrix} \mathbf{w}_{\ell} \\ b_{\ell} \end{bmatrix} = \frac{\zeta_{\ell'}}{c_{\ell,1}} \begin{bmatrix} \Phi_{\ell'}^{\top} \alpha_{\ell'} + \Phi_{\ell}^{\top} A_{\ell}^{-1} \hat{K}_{\ell,\ell'} \alpha_{\ell'} \\ \mathbf{1}_{\ell'}^{\top} \alpha_{\ell'} + \mathbf{1}_{\ell}^{\top} A_{\ell}^{-1} \hat{K}_{\ell,\ell'} \alpha_{\ell'} \end{bmatrix}, \quad (10.23)$$

entonces, la clasificación basada en ETWSVM no-lineal puede ser extraída de las ecuaciones (9.50), (10.22) y (10.23) como sigue:

$$t_* = \arg \min_{\ell \in \{+, -\}} \frac{|f_\ell(\mathbf{x}_*)|}{\|\mathbf{w}_\ell\|_{\mathcal{H}}}, \quad (10.24)$$

donde $\|\mathbf{w}_\ell\|_{\mathcal{H}}^2 = \langle \mathbf{w}_\ell, \mathbf{w}_\ell \rangle_{\mathcal{H}}$, y es calculado como:

$$\|\mathbf{w}_\ell\|_{\mathcal{H}}^2 = \frac{1}{c_{\ell,1}^2} \boldsymbol{\alpha}_{\ell'}^\top \left(\mathbf{K}_{\ell',\ell'} - 2\mathbf{K}_{\ell',\ell} \mathbf{A}_\ell^{-1} \hat{\mathbf{K}}_{\ell,\ell'} + \hat{\mathbf{K}}_{\ell',\ell} \mathbf{A}_\ell^{-1} \mathbf{K}_{\ell,\ell} \mathbf{A}_\ell^{-1} \hat{\mathbf{K}}_{\ell,\ell'} \right) \boldsymbol{\alpha}_{\ell'} \quad (10.25)$$

En resumen, la extensión no-lineal de nuestro enfoque ETWSVM explota los fundamentos de la representación kernel a través de un mapeo adecuado a un espacio de Hilbert con kernel reproductivo durante la optimización de los hiperplanos; mientras que ofrece una forma cerrada para el calculo del puntaje y las etiquetas de salida basadas en matrices kernel controladas por datos.

10.2. Aprendizaje de la función kernel en ETWSVM usando alineación centrada

Dentro del reconocimiento de patrones, en promedio, se prefiere utilizar el kernel Gaussiano por su capacidad de generalización, aproximación universal, y trazabilidad matemática. De ahí, que el siguiente función kernel surja en nuestro clasificador no-Lineal ETWSVM:

$$\kappa_{\Sigma}(\mathbf{x}_n, \mathbf{x}_{n'}) = \exp\left(-\frac{1}{2}(\mathbf{x}_n - \mathbf{x}_{n'})^\top \Sigma^{-1}(\mathbf{x}_n - \mathbf{x}_{n'})\right), \quad (10.26)$$

donde $\Sigma \in \mathbb{R}^{P \times P}$ es la matriz de covarianza del kernel Gaussiano. Con el fin de evitar problemas de inestabilidad con respecto a la inversa de Σ y de encontrar estructuras de datos discriminativos, usamos una estrategia de alineación de kernel centralizado (*Centered Kernel Alignment*–(CKA)) para estimar Σ . En particular, CKA aprovecha un ajuste controlado por datos de la matriz de covarianza del kernel Gaussiano al cuantificar la similitud entre la matriz \mathbf{K} (con elementos dados por la ecuación (10.26)) y un kernel calculado sobre las etiquetas de salida $\mathbf{Z} \in \mathbb{R}^{N \times N}$. Este último, $z_{n,n'} = \delta(t_n - t_{n'})$, siendo $\delta(\cdot, \cdot)$ la función delta. Luego, reescribimos la inversa de la matriz de covarianza como: $\Sigma^{-1} = \mathbf{E} \mathbf{E}^\top$, donde $\mathbf{E} \in \mathbb{R}^{P \times P'}$ ($P' \leq P$), y redefinimos la ecuación (10.26) explorando una proyección lineal, para que:

$$\kappa_{\Sigma}(\mathbf{x}_n, \mathbf{x}_{n'}; \mathbf{E}) = \exp\left(-\frac{1}{2}\|\mathbf{x}_n \mathbf{E} - \mathbf{x}_{n'} \mathbf{E}\|_2^2\right). \quad (10.27)$$

La función kernel en la ecuación (10.27) fomenta la alineación entre las matrices \mathbf{K} y \mathbf{Z} según la siguiente función de costo basada en CKA [60]:

$$\rho(\mathbf{K}(\mathbf{E}), \mathbf{Z}) = \left(\frac{\langle \widetilde{\mathbf{K}}(\mathbf{E}), \widetilde{\mathbf{Z}} \rangle_F}{\|\widetilde{\mathbf{K}}(\mathbf{E})\|_F \|\widetilde{\mathbf{Z}}\|_F} \right), \quad (10.28)$$

donde $\widetilde{\mathbf{K}} \in \mathbb{R}^{N \times N}$ representa una matriz kernel centralizado calculado como $\widetilde{\mathbf{K}} = \widetilde{\mathbf{I}} \mathbf{K} \widetilde{\mathbf{I}}$, $\widetilde{\mathbf{I}} = \mathbf{I} - \mathbf{1}\mathbf{1}^\top / N$; donde $\langle \cdot, \cdot \rangle_F$ y $\|\cdot\|_F$ denotan el producto interno y la norma de Frobenius basada en matrices, respectivamente. Además, $\widetilde{\mathbf{K}}(\mathbf{E})$ destaca la dependencia del kernel Gaussiano sobre la proyección lineal \mathbf{E} . En efecto, la función de costo 10.28 es usado para aprender la inversa de la matriz de covarianza Σ^{-1} desde \mathbf{E} como sigue [38]:

$$\mathbf{E}_* = \arg \max_{\mathbf{E}} \log(\rho(\mathbf{K}(\mathbf{E}), \mathbf{Z})), \quad (10.29)$$

donde la función log aumenta la convergencia del algoritmo. El problema de optimización en la ecuación (10.29) puede ser resuelta empleando el enfoque de gradiente descendiente [37]:

$$\begin{aligned} \nabla_E(\rho(\mathbf{K}(\mathbf{E}), \mathbf{Z})) = & -4\mathbf{X}((\mathbf{G} \circ \mathbf{K}(\mathbf{E})) - \dots \\ & \dots \text{diag}(\mathbf{1}^\top(\mathbf{G} \circ \mathbf{K}(\mathbf{E})))\mathbf{X}^\top \mathbf{E}, \end{aligned} \quad (10.30)$$

donde $\text{diag}(\cdot)$ y el \circ denota el operador diagonal y el producto de Hadamard, y $\mathbf{G} \in \mathbb{R}^{N \times N}$ es el gradiente de la función objetivo con respecto a $\mathbf{K}(\mathbf{E})$, calculada como:

$$\mathbf{G} = \frac{\tilde{\mathbf{Z}}}{\text{tr}(\mathbf{K}(\mathbf{E})\tilde{\mathbf{Z}})} - \frac{\tilde{\mathbf{K}}(\mathbf{E})}{\text{tr}(\mathbf{K}(\mathbf{E})\tilde{\mathbf{K}}(\mathbf{E}))}, \quad (10.31)$$

donde $\text{tr}(\cdot)$ representa el operador traza. Notablemente, nuestra estrategia basada en CKA brinda una similitud basada en Gauss que favorece la separabilidad de los datos a través de la concordancia entre la entrada y el kernel basado en las etiquetas. Por otra parte, se evita la inversión de la matriz de covarianza debido a la relación $\Sigma^{-1} = \mathbf{E}_* \mathbf{E}_*^\top$.

10.3. ETWSVM para problemas Multi-clase

Para los problemas de clasificación multi-clase, la base de datos desbalanceada contiene etiquetas de salidas $t_n \in \{1, 2, \dots, R\}$, donde $R \in \mathbb{N}$ es el número de clases en \mathbf{X} . Entonces, la matriz de muestras de la r -ésima clase $\mathbf{X}_r \in \mathbb{R}^{P \times N_r}$ es construida, donde $N = \sum_{r=1}^R N_r$ y $r \in \{1, 2, \dots, R\}$.

10.3.1. ETWSVM One-versus-Rest

El esquema *One-versus-Rest* (OvR) es bien conocido y es uno de los más populares enfoques para problemas multiclase a partir de clasificadores binarios. Este sería acoplado con nuestro ETWSVM para la construcción de R hiperplanos no-paralelos y por ende R QPPs independientes, uno por cada clase [25]. Formalmente, ETWSVM-OvR encuentra R hiperplanos de separación $f_r(\mathbf{x}) = \varphi(\mathbf{x}^\top \mathbf{w}_r) + b_r = 0$, como se muestra en la figura 10.1(a), tal que:

$$\begin{aligned} \min_{\mathbf{w}_r, b_r, \xi_{r'}} \quad & \frac{c_{1,r}}{2} (\|\mathbf{w}_r\|_{\mathcal{H}}^2 + b_r^2) + \frac{1}{2} \|\Phi_r^\top \mathbf{w}_r + \mathbf{1}_r b_r\|_{\mathcal{H}}^2 + c_{2,r'} \mathbf{1}_{r'}^\top \xi_{r'} \\ \text{s.t.} \quad & (\Phi_r^\top \mathbf{w}_r + \mathbf{1}_r b_r) \zeta_{r'} \geq \mathbf{1}_{r'} - \xi_{r'}, \\ & \xi_{r'} \geq \mathbf{0}, \quad \forall r \in \{1, \dots, R\}; \end{aligned} \quad (10.32)$$

donde $\Phi_r \in \mathbb{R}^{Q \times N_r}$ reúne las muestras mapeadas de la clase r mientras que $\Phi_{r'} \in \mathbb{R}^{Q \times N_{r'}}$ las muestras mapeadas de las demás clases. Adicional, la formulación dual de el problema de optimización de la ecuación (10.32) es llevado acabo a través de lema de la matriz inversa para resolver R QPPs como en la ecuación (10.20), fijando $\mathbf{K}_{r,r'} = \Phi_r^\top \Phi_{r'} \in \mathbb{R}^{N_r \times (N - N_r)}$. Consecutivamente, la etiqueta de salida de una nueva muestra, \mathbf{x}_* , puede ser inferida como:

$$t_* = \arg \min_{r \in \{1, \dots, R\}} \frac{|f_r(\mathbf{x}_*)|}{\|\mathbf{w}_r\|_{\mathcal{H}}}, \quad (10.33)$$

donde $f_r(\mathbf{x}_*)$ y $\|\mathbf{w}_r\|_{\mathcal{H}}$ con calculados como en las ecuaciones (10.22) y (10.25), respectivamente.

10.3.2. ETWSVM One-versus-One

Por otra parte, el esquema *One-versus-One* (OvO) acoplado con nuestro ETWSVM busca construir $R(R-1)/2$ clasificadores binarios [61], como se muestra en la figura 10.1(b). Explicitamente, para discriminar entre las clases r y r' , ETWSVM-OvO almacena las matrices $\Phi_r \in \mathbb{R}^{Q \times N_r}$ y $\Phi_{r'} \in \mathbb{R}^{Q \times N_{r'}}$ para solucionar el QPPs de la ecuación (10.32) solamente para las dos clases de interés. En la etapa de decisión, ETWSVM-OvO adopta la estrategia de “voto mayoritario”. Por ejemplo, si para una nueva muestra x_* el sub-clasificador entre las clases r -ésima y r' -ésima identifica la nueva etiqueta como r (evaluando la ecuación (10.33) únicamente para las dos clases bajo estudio), la votación para la r -ésima clase es aumentado en uno. Luego de evaluar todos los $R(R-1)/2$ sub-clasificadores, la etiqueta t_* es asignada por la clase con mayor número de votos [39].

Es importante mencionar que la estrategia de aprendizaje kernel presentada en la apartado 10.2 puede ser adoptado dentro de los ETWSVM-OvR y ETWSVM-OvO ya que los datos de entrada y los kernels basados en las etiquetas de salida (K y Z) no dependen de la tarea en estudio, es decir, si es para clasificación binaria o multi-clase.

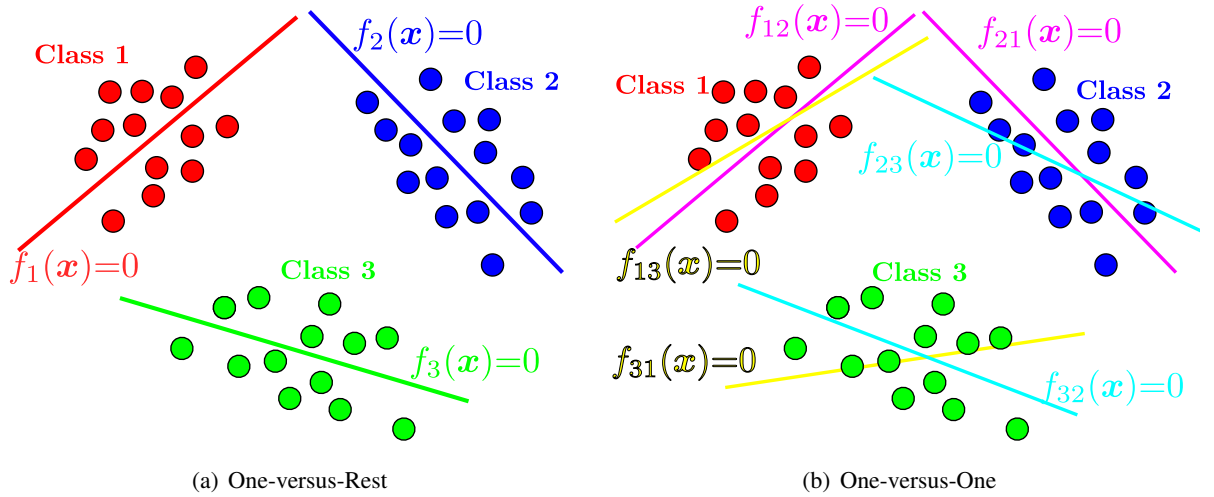


Figura 10.1.: Representación geométrica de las estrategias multiclase acopladas con el ETWSVM. a) corresponde al esquema OvR acoplado con el ETWSVM, donde se generan $R=3$ hiperplano, uno por cada clase. b) representa el esquema OvO acoplado con el ETWSVM, donde para $R=3$, se general tres sub-clasificadores (por cada uno dos hiperplanos).

11. Montaje experimental

En este capítulo explicaremos las bases de datos usadas para probar nuestro ETWSVM, tanto para el problema de clasificación binaria como multi-clase. También el proceso de entrenamiento de nuestra propuesta y de los métodos con los cuales nos comparamos para demostrar las ventajas de la primera. Por último, las medidas de rendimiento utilizadas para evaluar todos los clasificadores.

11.1. Bases de Datos

Para ilustrar el rendimiento de nuestro ETWSVM, consideramos un par de base de datos sintéticas en dos dimensiones ($P=2$) devotas a clasificación binarias desbalanceadas: i) Unas media lunas (ver figura 12.1(a)) la cual posee una estructura no-lineal, fijando $N_+=97$ y $N_-=227$, y ii) dos grupos de datos generados a partir de dos distribuciones Gaussianas (ver figura 12.2), lo cual nos permite imponer diferentes tasas de desbalance ($\{1:2, \dots, 1:10\}$ fijando $N_-=100$) y escenarios de traslape (variarnos la distancia Euclídea entre las medias de las clases desde el conjunto $[0, 10.14]$). Además, empleamos el públicamente disponible Repositorio de Aprendizaje de maquina UCI¹ para bases de datos de clasificación binaria como lo describen en [62, 63]. Con respecto al escenario multi-clase, también considerado el repositorio Keel² como es discutido en [25, 36]. Las bases de datos contempla un amplio rango de campos, tales como patologías, vehículos, ingeniería, informática biológica, financiera, etc., de diferentes tamaños y tasas de desbalance. Una pequeña descripción de las bases de datos de los repositorios UCI y Keel se describe en el cuadro 11.1.

Binary classsification (UCI repository)					Multi-class classsification (Keel repository)			
Name	IR	P	N	Minority class	Name	N	P	R
BankNote	0.8005	4	1372	1	Balance	625	4	3
Wisconsin	0.5938	30	569	Malignant	Contraceptive	1473	9	3
Ionosphere	0.5600	34	351	b	Dermatology	358	34	6
Pima-Indians	0.5360	8	768	1	Ecoli	336	7	8
Biodeg	0.5093	41	1055	RB	Glass	214	7	7
Iris	0.5000	4	150	Iris Virginica	Hayes-roth	160	4	3
Haberman	0.3600	4	306	2	New-thyroid	215	5	3
Transfusion	0.3123	4	748	1	Iris	150	4	3
Vehicle	0.3076	18	846	Van	Thyroid	7200	21	3
cmc	0.2921	9	1473	2:Long-term	Wine	178	13	3
Housing	0.0743	13	506	1 (CHAS)	Penbased	10992	16	10
Balance	0.0851	4	625	Balanced				

Tabla 11.1.: Repositorios UCI y Keel descritos para clasificación desbalanceada. $IR = N_+/N_-$: tasa de desbalance, P : # de características, N : # de muestras, y R : # de clases.

¹<https://archive.ics.uci.edu/ml/index.php>

²<https://sci2s.ugr.es/keel/category.php?cat=clas>

11.2. Entrenamiento de ETWSVM y métodos de comparación

11.2.1. Pre-procesamiento

Considerando la tarea de clasificación binaria, un ν -nearest neighbor (ν -NN), método de sub-muestreo, es llevado a cabo antes del aprendizaje del kernel de ETWSVM para eludir la influencia de los valores atípicos y para mitigar la carga computacional cuando inferimos la proyección lineal en la ecuación (10.27). Por consiguiente, un grafo local es definido para revelar la densidad intracase en el grupo mayoritario siguiendo una filosofía familiar como en [30]. Sea $U \in \mathbb{R}^{N_- \times N_-}$ una matriz adyacente definida con elementos:

$$u_{nn'} = \begin{cases} 1 & \mathbf{x}_n \in \Omega_\nu(\mathbf{x}_{n'}) \text{ y } \mathbf{x}_{n'} \in \Omega_\nu(\mathbf{x}_n) \\ 0 & \text{Otro caso,} \end{cases} \quad (11.1)$$

donde Ω_ν contiene los $\nu \in \mathbb{N}$ vecinos más cercanos de \mathbf{x}_n en X_- acorde a la distancia Euclidea ($\forall \mathbf{x}_n, \mathbf{x}_{n'} \in X_-$). Luego, es calculado el vector, $\vartheta \in \mathbb{R}^{N_-}$, de coeficientes de sub-muestreo como: $\vartheta_n = \sum_{\mathbf{x}_{n'} \in X_-} u_{nn'}$. La matriz de la clase mayoritaria puede ser reducida por el número de muestras de la clase minoritaria, N_+ con los valores más altos del coeficiente de sub-muestreo en ϑ . En consecuencia, las muestras provenientes de regiones de alta densidad son más propensos a ser seleccionados que las muestras provenientes de regiones de baja densidad, por ejemplo, se descartan los valores atípicos. Del mismo modo, para la tarea multi-clase, la tarea del sub-muestreo es llevada a cabo para eliminar datos de entrada forzando los cada uno de los $R - 1$ clases (excluyendo la clase con más bajo número de muestras) a contener el siguiente número de instancias: $\hat{N} = \min(\min_{r'} R N_r, N_r)$. El factor R alienta la preservación de la estructura de los datos de entrada respecto al número de clases proporcionadas. Para concretar la prueba, nosotros experimentalmente fijamos $\nu=8$ inspirado en una grilla de búsqueda discutida por los autores en [30].

11.2.2. Aprendizaje de la función kernel

La función kernel es obtenida a través del enfoque basado en CKA discutido en el apartado 10.2, configurando los siguientes parámetros de optimización: La inicialización de la proyección lineal dentro del iterativo gradiente descendiente es ajustada acorde al bien conocido enfoque *Principal Component Analysis* (PCA), la tolerancia del gradiente descendiente es fijado en 10^{-6} , el máximo número de iteraciones es limitado empíricamente a 300, y el número de dimensiones relevantes (P^h) es ajustada para retener el 95 % de la varianza explicada [38]. En aras de simplicidad, consideramos ETWSVM* como el emparejamiento entre el preprocesamiento basado en ν -NN (apartado 11.2.1), el aprendizaje de la función kernel basado en CKA (apartado 11.2.2), y la discriminación basada en el ETWSVM (ver figura 11.1). Asimismo, ETWSVM*-OvO y ETWSVM*-OvR sostienen un multi-clase extensión de ETWSVM*.

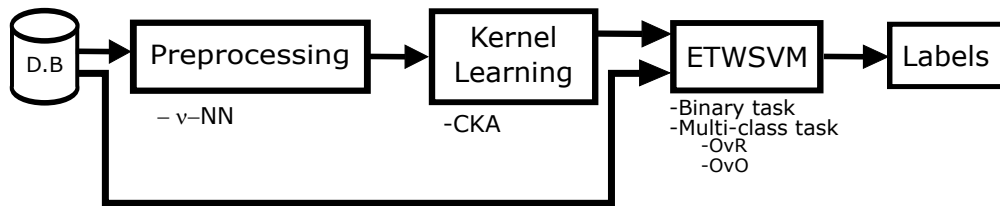


Figura 11.1.: Disposición del ETWSVM*. Un submuestreo basado en ν -NN es aplicado antes de la etapa del aprendizaje de la función kernel. Entonces, la clasificación binaria o multiclase es realizado para resolver el problema de optimización de ETWSVM.

11.2.3. Evaluación y sintonización de ETWSVM

Para la sintonización, de los parámetros libres de los clasificadores, y la evaluación de los mismos, tanto binarios como multiclase, usamos un esquema de validación cruzada anidada 10-*folds*. Para el escenario binario, usamos las medidas de rendimiento *Accuracy*-(*Acc*), la *Geometric Mean*-(*GM*), y la *F-measure*-(*FM*), las cuales son calculadas como sigue:

$$\begin{aligned} Acc &= \frac{T_P + T_N}{N_{test}} \\ GM &= \sqrt{Sen \times Spec} \\ FM &= \frac{2Pre \times Rec}{Pre + Rec}, \end{aligned} \quad (11.2)$$

donde $T_P \in \mathbb{N}$ es el número de muestras positivas que son estimadas como positivas, $T_N \in \mathbb{N}$ es el número de muestras negativas que son estimadas como negativas, y $N_{test} \in \mathbb{N}$ es el número de instancias en el conjunto de prueba. $Sen = T_P / (T_P + F_N)$ y $Spec = T_N / (T_N + F_P)$ son las medidas de sensibilidad y especificidad, respectivamente; donde $F_N \in \mathbb{N}$ es el número de muestras positivas que son clasificadas como negativas y $F_P \in \mathbb{N}$ es el número de muestras negativas que son clasificadas como positivas. Por otra parte, $Pre = T_P / (T_P + F_P)$ y $Rec = Sen$ son los valores de las medidas de precisión y recuperación, respectivamente. En general, el *Acc* es la medida de rendimiento, en clasificación, más popular; sin embargo, esta no es indicada para calcular el rendimiento en problemas de clasificación desbalanceada, un ejemplo es supongamos una base de datos la cual tenga una tasa de desbalance es de 9 : 1, si el clasificador sesga por completo a la clase minoritaria y clasifica bien todas las muestras provenientes de la clase mayoritaria su $Acc = 90\%$. En consecuencia, la evaluación basada en *GM* y *FM* son también adoptadas, ya que tienen menos probabilidades de sufrir el problema de desbalance, dado que tienen en cuenta la distribución de clase [8].

Adicional, para evaluar el escenario multi-clase, contemplamos las medidas de rendimiento *Mean Accuracy*-(\overline{Acc}) y *F-measure* $_{\mu}$ -(FM_{μ}). El primero captura la efectividad promedio del clasificador junto con las R clases mientras el último favorece la capacidad predictiva considerando las etiquetas positivas [64]. Estas medidas son calculadas como sigue:

$$\begin{aligned} \overline{Acc} &= \frac{1}{R} \sum_{r=1}^R \frac{T_{P_r} + T_{N_r}}{N_{test}} \\ FM_{\mu} &= \frac{2 \overline{Pre} \times \overline{Rec}}{R \overline{Pre} + \overline{Rec}}, \end{aligned} \quad (11.3)$$

donde $\overline{Pre} = \frac{1}{R} \sum_{r=1}^R Pre_r$ y $\overline{Rec} = \frac{1}{R} \sum_{r=1}^R Rec_r$. El subíndice r indica que la cantidad correspondiente se computará en la clase r -ésima.

Los parámetros de regularización son buscados del conjunto $\{2^{-7}, 2^{-5}, \dots, 2^5, 2^7\}$ con respecto a la evaluación basada en *FM*. No obstante, para la reducción de la complejidad computacional, en nuestros experimentos, fijamos $c_{1,\ell} = c_{1,\ell'}$ y $c_{2,\ell} = c_{2,\ell'}$. Nótese que el ETWSVM sencillo incluye un kernel Gaussiano con covarianza isotrópica $\Sigma = \sigma^2 \mathbf{I}$, donde el valor del ancho bando $\sigma^2 \in \mathbb{R}^+$ es buscado dentro del rango $\{0.1\sigma_0, 0.2\sigma_0, \dots, \sigma_0\}$ con respecto al valor de *FM*, $\sigma_0 = \text{med}(\|\mathbf{x}_n - \mathbf{x}_{n'}\|_2)$ y $\text{med}(\cdot)$ indica el operador mediana (no se aplica la etapa de preprocesamiento). Todos los QPPs relacionados con nuestro ETWSVM son resueltos utilizando el algoritmo de punto-interior-convexo del Toolbox de Optimización de MATLAB, fijando la tolerancia de el método en 10^{-12} , la terminación en el primer modo óptimo y las tolerancias de

violación de las restricciones en 10^{-8} , y el máximo número de iteraciones en 200.

11.2.4. Métodos de comparación

Junto al ETWSVM sencillo (el cual incluye un kernel Gaussiano con covarianza isotrópica) y el ETWSVM* (el cual contiene el preprocesamiento basado en ν -NN y el aprendizaje de la función kernel basado en CKA), también contemplamos el bien conocido método de sobre-muestreo SMOTE en nuestra propuesta. Entonces, nos referimos a ETWSVM** como la implementación del algoritmo SMOTE para mejorar la tasa de desbalance; entonces, la base de datos con sobre-muestreo es usada para el aprendizaje de kernel basado en CKA y la construcción de la frontera de decisión. El número de vecinos más cercanos dentro del SMOTE es fijado desde el conjunto $\{1, 3, 6, 9, 12\}$ como es discutido en [18]. Aquí, SMOTE es únicamente usado con clasificadores binarios, debido a su creciente costo computacional [7, 65]. Sin embargo, la evaluación del ETWSVM** en problemas binarios nos dará información sobre las ventajas y desventajas del algoritmo SMOTE en comparación al ETWSVM y ETWSVM*. La implementación en MATLAB de nuestros ETWSVM, ETWSVM* y ETWSVM** (incluyendo ambas opciones, binaria y multi-clase) están públicamente disponibles³.

Para realizar una buena comparación, los siguientes métodos del estado del arte dedicados a la clasificación binaria son examinados: i) SVM [66], ii) WSVM [67], iii) TBSVM [59], iv) WLTSVM [30], y v) SVM_{SMOTE} [68]. El toolbox de estadística y aprendizaje de máquina de MATLAB son usados para implementar los clasificadores SVM y WSVM. Los scripts de WLTSVM son posteados por sus autores en [30]⁴. Asimismo, los algoritmos TBSVM y SMOTE son implementados como se menciona en [18, 59]. La versión no-lineal prueba ser más competitiva que la alternativa lineal para todos los métodos proporcionados en el estado del arte, y garantiza una comparación justa con nuestros algoritmos basados en ETWSVM, fijando un kernel Gaussiano isotrópico. El ancho de banda, σ , es buscado como en el ETWSVM sencillo. Adicional, los siguientes enfoques multi-clase son estudiados: SVM-OvR, SVM-OvO, TBSVM-OvR, y TBSVM-OvO [39]. De nuevo, una función kernel Gaussiana isotrópica es utilizada. Todos los experimentos fueron realizados en el entorno de MATLAB 9.1 en una computadora con un procesador i5 de 4^{ta} generación (2.0GHz) y 6GB de memoria RAM.

³https://github.com/cralji/ETWSVM_PR.git

⁴<http://www.optimal-group.org/Resource/LTSVM.html>

Parte III.

Resultados, Conclusiones y Trabajos Futuros

12. Resultados y Discusión

12.1. Resultados en datos sintéticos

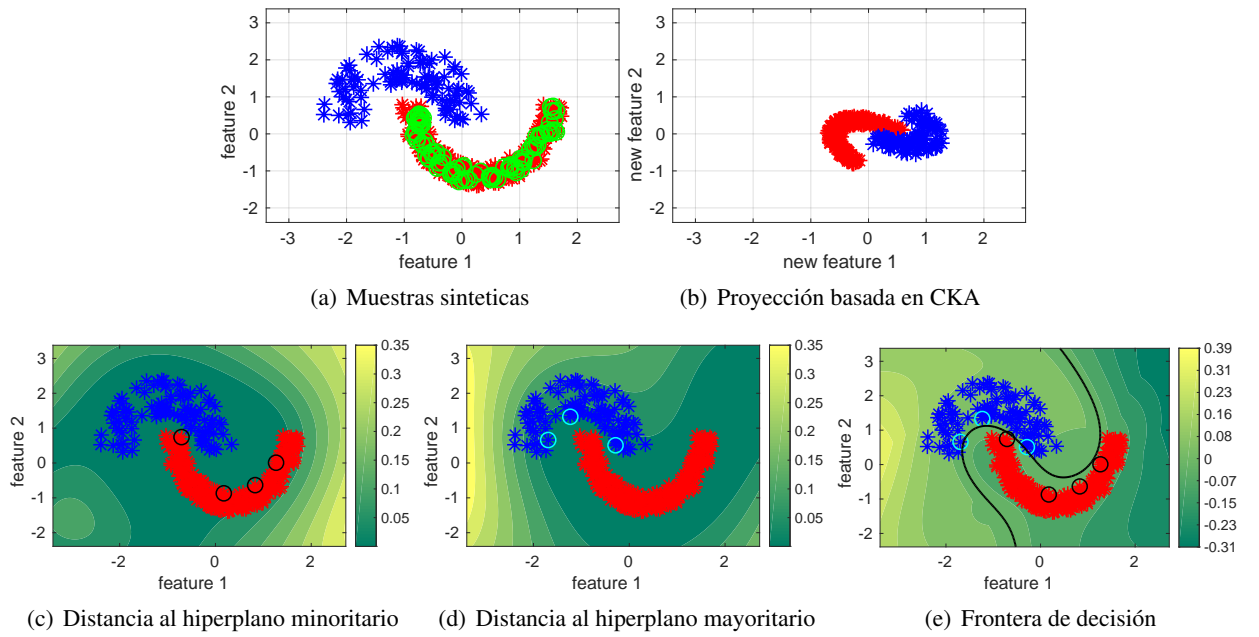


Figura 12.1.: Resultados de la clasificación del ETWSVM* sobre la base de datos media-luna. (a) muestra los datos sintéticos de media-luna, las muestras azules provienen de la clase minoritaria mientras las muestras rojas de la clase mayoritaria; las instancias encerradas en círculos verdes son seleccionadas por el submuestreo basado en ν -NN. (b) muestra la proyección basada en CKA fijando $P'=2$. (c) y (d) muestra la distancia entre los datos de entrada y los hiperplanos minoritario y mayoritario, respectivamente; donde las muestras encerradas en negro y en cian representa los vectores de soporte de los hiperplanos minoritario y mayoritario, respectivamente. (e) muestra la frontera de decisión y la mínima distancia entre (c) y (d) (ver función de decisión del ETWSVM en ecuación (10.24)).

Primero, un ejemplo representativo es llevado a cabo sobre la base de datos de media-luna para inspeccionar visualmente el rendimiento del ETWSVM*. Inicialmente, nosotros sólo describimos los resultados del ETWSVM* para dilucidar las virtudes de la representación respecto al mapeo no-lineal para un sostenible RKHS. En la figura 12.1(a) se muestra los datos de entrada y las muestras seleccionadas de la clase mayoritaria a través de la etapa del preprocesamiento basado en el ν -NN. Como se puede ver la etapa del ν -NN captura las muestras relevantes codificando la estructura de luna. Junto a, la proyección basada en

CKA, en este experimento lo fijamos $P'=2$ para visualizar el propósito, las metas para encontrar un compromiso entre “estiramiento” y “aplastamiento” para revelar información relevante (ver figura 12.1(b)). El espacio proyectado, el cual codifica dentro de la Gaussiana multi-variada (ver ecuación (10.27)), compensando la diferencia entre las densidades de las clases mayoritaria y minoritaria colapsando las instancias adecuadamente. Notablemente, las “nuevas nubes de puntos” concentrando la estructura fundamental de cada grupo. Por lo tanto, la distancia a los hiperplanos minoritarios y mayoritarios son mostrados en las figuras 12.1(c) y 12.1(d) exhibiendo como ETWSVM* codificando el funcionamiento de la media-luna. En efecto, ETWSVM* sólo requiere siete vectores de soporte ($\alpha > 0$) para construir una frontera suave como se ve en la figura 12.1(e).

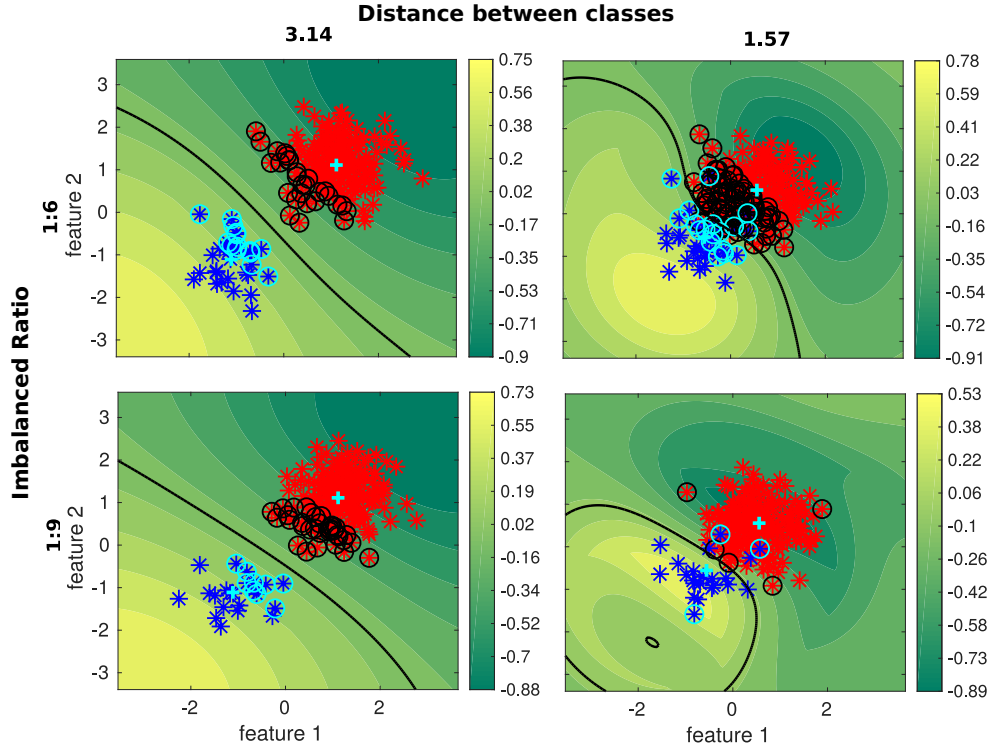


Figura 12.2.: Inspección visual de los resultados sobre bases de datos sintéticas generada a partir de dos distribuciones Gaussianas multivariadas (clasificación basada en ETWSVM*). Diferentes escenarios de desbalance y traslape entre las clases son probadas. Las muestras azules provienen de la clase minoritaria mientras que las rojas pertenecen a la mayoritaria. Las muestras encerradas en negro y cían representan los vectores de soporte de las clase minoritaria y mayoritaria, respectivamente, y la cruz la media de las distribuciones. La frontera de decisión es representada con una línea negra y la mínima distancia a los hiperplanos es coloreada en el fondo.

Segundo, la figura 12.2 presenta las fronteras de decisión del ETWSVM* para las bases de datos generadas a partir de dos Gaussianas multivariadas considerando varios escenarios de desbalance y traslape entre las clases. Como se puede apreciar, una tasa alta de desbalance influencia a la frontera de decisión en estar más cerrada a la clase minoritaria; también, un alto traslape entre clases, causa un incremento en los vectores de soporte de los hiperplanos. A demás, las figuras 12.3 y 12.4 muestran el porcentaje de GM y FM resultado de la clasificación sobre bases de datos generadas a partir de dos distribuciones Gaussianas en función de la tasa de desbalance, el traslape de las clases, y los valores de los parámetros de regularización. En efecto, tales parámetros fueron variados desde el conjunto $c_{1,\ell}, c_{2,\ell} \in \{10^{-4}, 10^{-1}, 1, 10^1, 10^2\}$, y fijamos la función

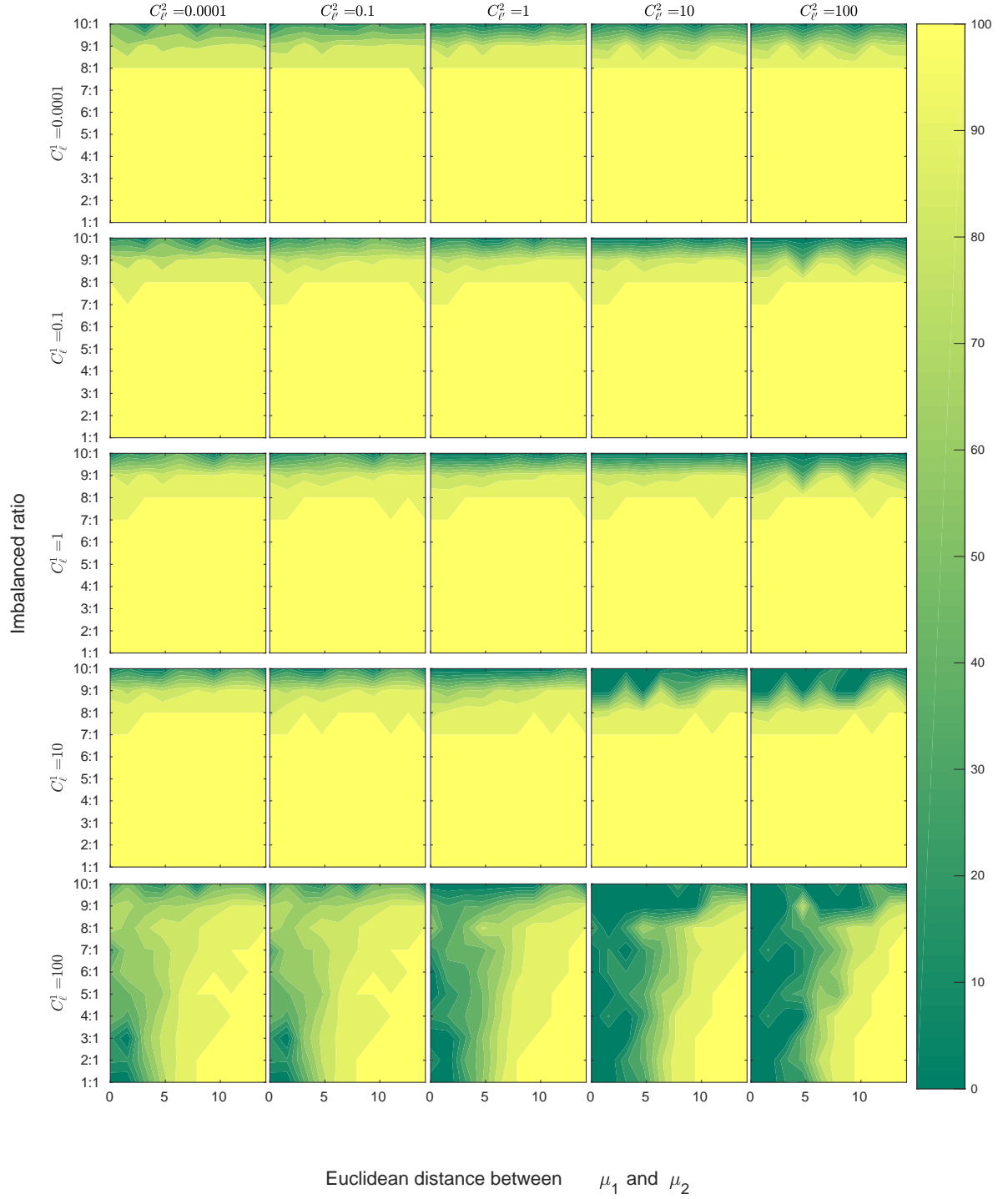


Figura 12.3.: Análisis de los parámetros de regularización (resultados de clasificación del ETWSVM*). La evaluación cualitativa basada en el GM se representa variando la tasa de desbalance y la distancia entre las medias de cada clase. Dentro de cada gráfica un par de valores de $c_{1,\ell}$ y $c_{2,\ell}$ es fijado para solucionar los QPPs en la ecuación (10.20). Esta medida de rendimiento se calcula variando de abajo hacia arriba la relación de desbalance y de izquierda a derecha la distancia entre las clases.

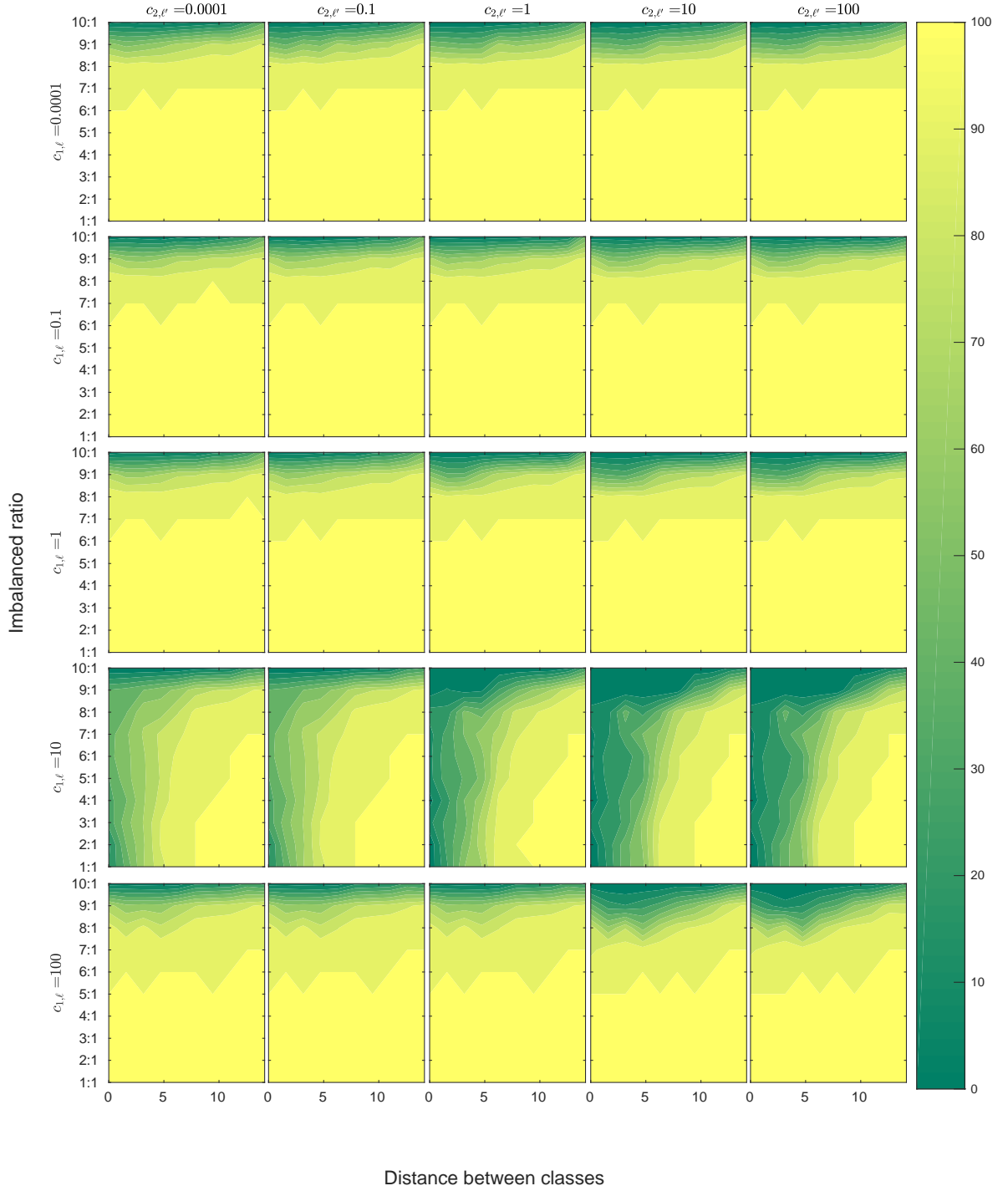


Figura 12.4.: Análisis de los parámetros de regularización (resultados de clasificación del ETWSVM*). La evaluación cualitativa basada en el FM se representa variando la tasa de desbalance y la distancia entre las medias de cada clase. Dentro de cada gráfica un par de valores de $c_{1,\ell}$ y $c_{2,\ell'}$ es fijado para solucionar los QPPs en la ecuación (10.20). Esta medida de rendimiento se calcula variando de abajo hacia arriba la relación de desbalance y de izquierda a derecha la distancia entre las clases.

kernel a través de la proyección basada en CKA. Cómo es claro, el parámetro $c_{1,\ell}$ es más sensible a la tasa de desbalance, y el clasificador refleja mejores rendimientos para valores bajos de $c_{1,\ell}$. Lo anterior puede ser explicado por la influencia de tal parámetro en la matriz inversa de la matriz Hessiana (ver ecuación (10.19)). Por el contrario, el parámetro $c_{2,\ell'}$ permite compensar el traslape de las clases por la penalización de la mala-clasificación de las muestras provenientes de la clase $\zeta_{\ell'}$ (ver ecuación (10.2)). De ahí, después de la inspección visual en los resultados del FM en la figura 12.4 podemos fijar $c_{1,\ell}=10^{-3}$ y $c_{2,\ell'}=1$ para obtener una compensación entre el rendimiento del clasificador y el costo computacional llevado a cabo por la sintonización de estos parámetros, cediendo a las siguientes extensiones del ETWSVM: ETWSVM[†], ETWSVM^{†*} y ETWSVM^{†**} († presenta los valores de los parámetros de regularización fijados).

12.2. Resultados en bases de datos reales: clasificación binaria

Para ilustrar lo propuesto, aplicamos el bien conocido algoritmo *t-student stochastic neighbor embedding* (*t-SNE*) para incrustar la base de datos Vehicle dentro de un espacio de características de dos dimensiones (2D) [69]. Aquí, nuestro objetivo es visualizar la capacidad de generalización de nuestro algoritmo respecto a la compleja frontera de decisión en el problema de clasificación en el mundo real. La figura 12.5 muestra la frontera de decisión de: SVM, TBSVM, ETWSVM, y ETWSVM*, ETWSVM[†], y ETWSVM^{†*} en el espacio 2D embebido por el t-SNE. Como se puede observar, el SVM no es capaz de codificar la estructura fundamental de los datos y esta completamente sesgado hacia la clase mayoritaria. El TBSVM construye una frontera de decisión que encierra la clase minoritaria, lo cual puede ser producto de sobre-aprendizaje sobre esta clase. Por otra parte, el ETWSVM y ETWSVM[†], los cuales incluyen un mapeo adecuado a RKHS, pueden codificar las estructuras de la clase mayoritaria y minoritaria pero a expensas de una frontera de decisión compleja. Adicional, el ETWSVM* y ETWSVM^{†*}, los cuales incluyen el aprendizaje del kernel por medio de una proyección lineal basada en CKA, puede revelar la información discriminativa manteniendo una frontera de decisión suave que garantiza una clasificación confiable con capacidad de generalización.

Sucesivamente, las tablas 12.1 y 12.2 presentan el rendimiento de clasificación y tiempo de entrenamiento para las bases de datos del repositorio UCI. Notamos que nuestro ETWSVM y sus variaciones (ETWSVM*, ETWSVM**, ETWSVM[†], ETWSVM^{†*} y ETWSVM^{†**}); muestran resultados aceptables en comparación a los clasificadores SVM y WSVM sobre las bases de datos Haberman, Housing, Vehicle, Transfusion, Ionosphere, Balance, Biodeg, e Iris. Para el resto de bases de datos nuestros clasificadores superan a SVM y WSVM. Con respecto al algoritmo SVM_{SMOTE}, los clasificadores basados en ETWSVM son comparables con respecto al desempeño de la clasificación para todas las bases de datos estudiadas. Sin embargo, la clasificación por medio del SVM_{SMOTE} requiere un alto costo computacional debido a su etapa de sobre-muestreo (ver la tabla 12.2). En particular, nuestros métodos basados en ETWSVM pueden lidiar con la desafiante tasa de desbalance, incluso en ausencia de una etapa de re-muestreo para la clasificación. Además, los resultados del WLTSVM son comparables con nuestras propuestas en las bases de datos Haberman, Prima-Indians, BankNote e Iris. No obstante, para el resto de problemas de clasificación binaria, nuestros clasificadores exceden al rendimiento del WLTSVM. Ahora, para el algoritmo TBSVM, se pueden observar altos valores de Acc, pero tiende a obtener resultados sesgados, por ejemplo, viendo el desempeño con respecto al GM y el FM sobre las bases de datos Housing, Transfusion, y Balance.

En promedio, nuestros planteamientos logran los más altos valores de Acc, GM, y FM en promedio, y ganan en 10 de 12 bases de datos para todas las medidas de rendimiento calculadas. Notablemente, el ETWSVM*, el cual incluye el preprocesamiento basado en ν -NN y el aprendizaje de la función kernel por medio de la proyección lineal basada en CKA, consigue los siguientes resultados en promedio: 86 % Acc, 79 % GM, y 69 %, los cuales son estadísticamente similares a los resultados de ETWSVM y ETWSVM**. Además,

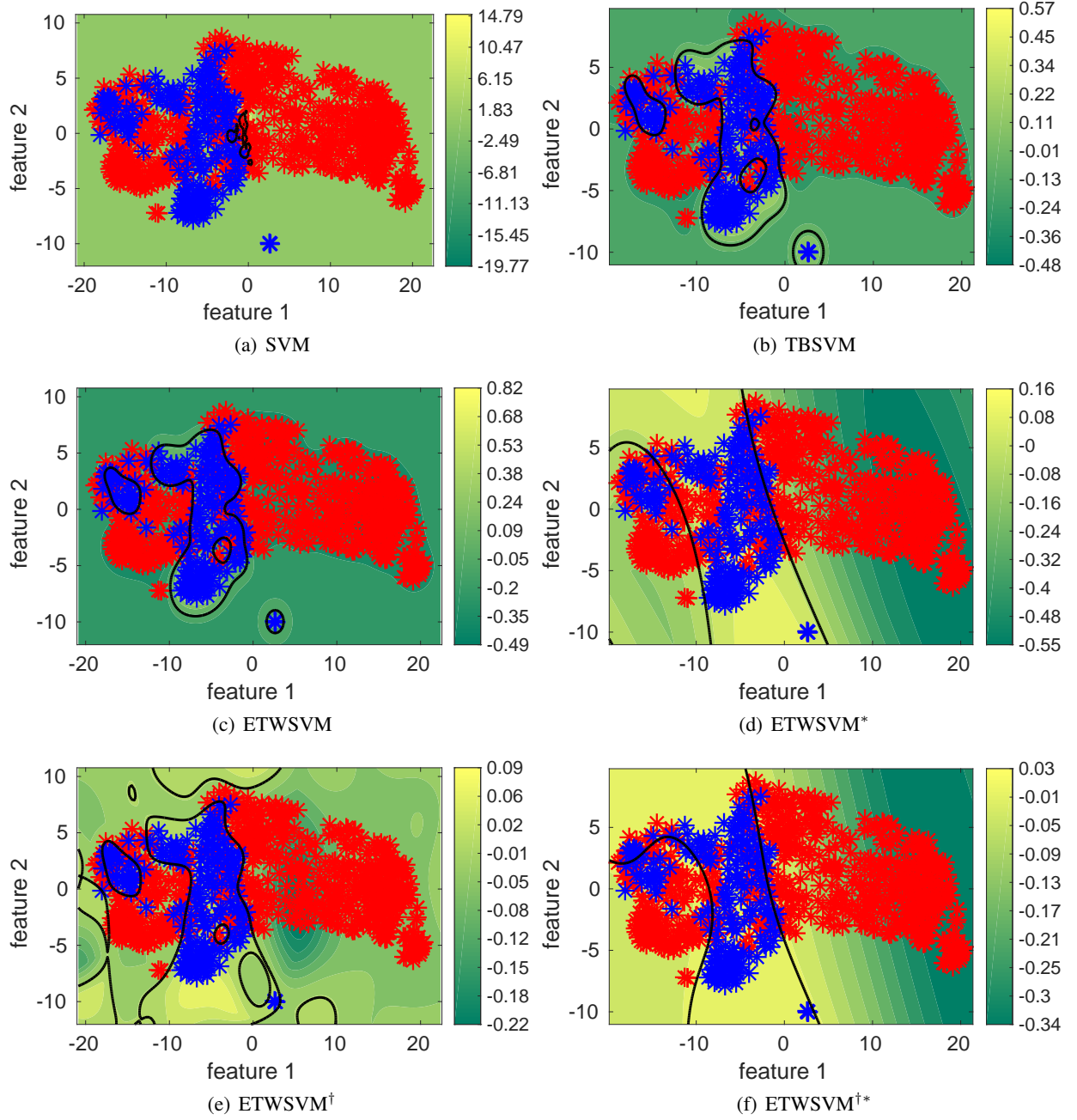


Figura 12.5.: Frontera de decisión de la base de datos Vehicle después de la embebimiento a un espacio 2D a través del t-SNE. Las muestras azules provienen de la clase minoritaria mientras que las rojas pertenecen a la clase mayoritaria. Las fronteras de decisión son mostradas en negro, y el puntaje de la clasificación (o la mínima distancia entre a los hiperplanos) son coloreados en el fondo.

después establecimos los valores de los parámetros de regularización en ETWSVM* como se explica en el apartado 12.1, ETWSVM^{†*}, el desempeño muestra una ligera disminución (87 % Acc, 72 % GM, y 64 % FM) pero se obtiene una ganancia significativa con respecto al tiempo computacional (ver tabla 12.2). De

	<i>SV M</i>	<i>WSVM</i>	<i>SVM_{SMOTE}</i>	<i>WLTSVM</i>	<i>TBSVM</i>	<i>ETWSVM</i>	<i>ETWSVM*</i>	<i>ETWSVM**</i>	<i>ETWSVM[†]</i>	<i>ETWSVM^{†*}</i>	<i>ETWSVM^{†**}</i>
Dataset	Acc GM FM	Acc GM FM	Acc GM FM	Acc GM FM	Acc GM FM	Acc GM FM	Acc GM FM	Acc GM FM	Acc GM FM	Acc GM FM	Acc GM FM
BankNote	88, 3±3, 7	92, 3±13, 0	97, 5±1, 9	91, 8±1, 4	91, 9±2, 2	97, 5±1, 4	96, 7±1, 7	96, 2±1, 7	97, 6±1, 1	96, 7±1, 4	96, 4±1, 3
	88, 7±3, 6	86, 6±30, 5	97, 4±2, 1	92, 2±1, 3	91, 9±2, 1	97, 5±1, 4	96, 7±1, 7	96, 3±1, 6	97, 6±1, 1	96, 8±1, 3	96, 6±1, 2
	87, 8±3, 6	86, 4±30, 4	97, 1±2, 2	91, 3±1, 4	43, 9±30, 4	97, 2±1, 6	96, 3±1, 9	95, 8±1, 8	97, 3±1, 3	96, 4±1, 4	96, 1±1, 4
Wisconsin	51, 3±5, 8	97, 2±1, 5	97, 2±2, 1	80, 8±5, 6	93, 0±2, 9	98, 2±2, 2	97, 7±2, 0	98, 1±1, 9	97, 0±2, 8	96, 7±3, 4	97, 4±2, 4
	46, 6±9, 2	97, 0±1, 7	96, 8±2, 3	81, 5±5, 5	90, 4±4, 1	97, 9±2, 3	97, 2±2, 5	97, 7±2, 5	96, 4±2, 9	96, 1±4, 3	97, 0±2, 8
	60, 6±3, 0	96, 2±2, 0	96, 2±2, 8	77, 0±6, 1	54, 3±0, 5	97, 6±2, 9	96, 9±2, 8	97, 4±2, 7	96, 0±3, 7	95, 3±4, 8	96, 4±3, 3
Inosphere	94, 6±2, 5	94, 0±2, 5	93, 7±5, 0	82, 4±6, 8	92, 3±5, 9	94, 6±5, 0	93, 7±5, 7	93, 7±3, 7	94, 0±3, 2	88, 7±5, 0	88, 3±5, 3
	93, 6±3, 1	93, 1±3, 3	92, 2±6, 7	82, 5±7, 6	93, 2±5, 5	93, 5±6, 1	91, 8±8, 4	92, 4±4, 5	92, 7±4, 0	84, 3±8, 2	83, 5±7, 7
	92, 3±3, 5	91, 5±3, 6	90, 8±7, 8	78, 3±8, 2	51, 3±5, 6	92, 2±7, 4	90, 4±9, 7	90, 9±5, 4	91, 4±4, 8	82, 0±9, 5	81, 1±9, 3
Prima-Indians	25, 3±4, 7	25, 0±4, 3	72, 8±5, 0	65, 4±7, 0	74, 4±3, 6	73, 6±3, 7	75, 5±3, 5	74, 0±5, 4	72, 5±5, 1	71, 5±2, 9	73, 8±4, 8
	24, 5±5, 6	24, 5±5, 0	72, 7±4, 6	68, 3±7, 1	71, 6±4, 1	72, 3±3, 7	73, 9±3, 1	73, 3±5, 4	72, 7±5, 7	65, 3±5, 2	69, 0±5, 9
	30, 3±6, 1	30, 3±5, 0	65, 4±5, 2	64, 7±5, 5	51, 7±0, 4	64, 8±4, 4	66, 6±3, 7	65, 9±6, 2	65, 3±6, 4	56, 1±6, 3	60, 7±7, 6
Biodeg	87, 0±2, 9	87, 3±3, 0	87, 8±2, 2	77, 8±2, 2	84, 1±5, 3	87, 4±1, 5	86, 6±1, 4	88, 1±2, 0	88, 7±1, 7	86, 4±4, 3	86, 8±3, 6
	86, 6±2, 8	86, 5±3, 0	86, 8±2, 7	75, 2±2, 2	84, 3±4, 6	85, 6±2, 7	84, 4±2, 6	86, 5±2, 6	86, 8±2, 6	85, 0±4, 8	85, 2±3, 9
	81, 7±3, 7	81, 8±3, 9	82, 3±3, 3	67, 7±2, 6	54, 8±7, 2	81, 2±2, 6	79, 8±2, 7	82, 3±2, 9	83, 0±2, 9	80, 2±6, 2	80, 6±5, 1
Iris	86, 0±9, 1	99, 3±2, 1	99, 3±2, 1	99, 3±2, 1	99, 3±2, 1	99, 3±2, 1	100, 0±0, 0	100, 0±0, 0	99, 3±2, 1	100, 0±0, 0	100, 0±0, 0
	73, 9±19, 3	98, 9±3, 3	98, 9±3, 3	98, 9±3, 3	98, 9±3, 3	98, 9±3, 3	100, 0±0, 0	100, 0±0, 0	98, 9±3, 3	100, 0±0, 0	100, 0±0, 0
	69, 8±23, 7	98, 9±3, 5	98, 9±3, 5	98, 9±3, 5	50, 0±0, 0	98, 9±3, 5	100, 0±0, 0	100, 0±0, 0	98, 9±3, 5	100, 0±0, 0	100, 0±0, 0
Haberman	62, 8±8, 5	71, 5±5, 0	61, 8±5, 2	59, 5±16, 3	73, 9±3, 9	73, 2±6, 2	68, 3±7, 8	70, 6±9, 0	61, 7±14, 8	71, 9±5, 3	69, 3±8, 3
	53, 5±12, 1	54, 5±9, 5	56, 2±12, 5	38, 3±22, 9	32, 0±24, 9	63, 8±10, 0	58, 5±7, 8	57, 0±14, 0	62, 5±13, 0	41, 0±10, 1	53, 5±10, 7
	37, 5±12, 9	50, 1±11, 4	40, 7±13, 5	28, 1±18, 1	41, 9±1, 2	50, 1±12, 0	44, 0±10, 0	42, 6±16, 2	49, 4±13, 8	26, 5±11, 1	38, 7±11, 7
Transfusion	69, 1±4, 1	68, 0±5, 0	62, 6±10, 3	36, 0±31, 2	23, 8±0, 4	76, 0±5, 7	76, 2±4, 6	73, 0±7, 2	75, 4±6, 2	79, 4±3, 4	74, 2±4, 9
	64, 8±5, 9	59, 7±22, 0	51, 9±28, 2	39, 2±33, 9	0, 0±0, 0	64, 6±7, 4	65, 7±7, 3	68, 3±7, 1	66, 1±6, 7	54, 7±8, 8	66, 4±8, 3
	47, 4±6, 2	44, 0±17, 0	39, 2±21, 6	49, 0±2, 9	38, 4±0, 6	50, 0±9, 6	51, 0±8, 8	52, 7±8, 5	51, 2±9, 6	42, 4±1, 1	50, 5±9, 3
Vehicle	97, 5±1, 4	97, 9±1, 3	97, 6±1, 5	82, 4±4, 6	85, 7±7, 0	97, 9±2, 1	95, 2±2, 9	96, 1±2, 4	98, 5±1, 4	96, 3±1, 3	96, 5±2, 4
	97, 7±2, 0	97, 7±1, 8	96, 9±2, 4	82, 2±4, 0	64, 8±44, 8	96, 8±3, 0	94, 0±5, 5	96, 5±2, 9	97, 9±2, 1	95, 3±2, 7	96, 4±2, 7
	94, 9±3, 0	95, 6±2, 7	95, 0±3, 1	69, 1±6, 4	38, 1±0, 5	95, 5±4, 2	89, 9±6, 5	92, 3±4, 7	96, 7±2, 9	92, 3±2, 8	92, 8±4, 8
cmc	26, 8±4, 7	33, 2±3, 2	67, 8±6, 2	31, 9±22, 4	64, 6±6, 3	68, 0±4, 1	66, 2±4, 2	65, 9±3, 3	67, 6±3, 5	75, 6±2, 7	68, 4±4, 0
	19, 9±17, 5	31, 7±4, 1	58, 8±21, 1	31, 7±22, 6	66, 0±5, 2	66, 6±3, 4	65, 2±3, 5	65, 6±3, 3	66, 8±5, 4	41, 5±7, 5	67, 1±5, 4
	21, 8±18, 9	44, 0±4, 3	41, 9±15, 5	28, 1±7, 3	36, 8±0, 6	47, 9±3, 5	46, 1±4, 0	46, 5±3, 7	47, 8±5, 9	26, 0±8, 0	48, 3±6, 0
Housing	91, 9±2, 7	88, 6±5, 3	90, 3±3, 7	44, 9±6, 6	93, 1±1, 0	89, 7±2, 7	88, 4±5, 9	83, 4±3, 8	89, 2±4, 1	91, 5±2, 3	86, 0±4, 0
	63, 4±16, 7	67, 7±27, 6	66, 3±25, 1	46, 2±18, 2	0, 0±0, 0	47, 8±29, 0	41, 6±38, 4	72, 3±11, 5	43, 7±39, 6	38, 7±27, 3	62, 0±26, 0
	42, 0±14, 9	41, 1±19, 4	42, 8±21, 2	13, 7±3, 7	12, 9±1, 7	27, 8±19, 9	27, 3±32, 3	34, 2±10, 2	26, 8±25, 9	24, 7±18, 1	31, 8±16, 3
Balance	90, 7±3, 3	90, 7±3, 2	83, 8±6, 0	49, 5±20, 7	92, 2±0, 5	83, 2±5, 9	84, 3±4, 7	83, 2±6, 7	84, 2±5, 0	93, 4±1, 8	85, 8±4, 6
	81, 1±18, 9	57, 3±11, 2	69, 3±13, 4	52, 5±16, 5	0, 0±0, 0	75, 8±13, 3	78, 4±13, 0	69, 8±21, 0	78, 7±17, 2	63, 5±13, 0	82, 2±10, 8
	54, 7±19, 4	59, 0±13, 2	37, 2±13, 6	18, 5±8, 8	11, 0±5, 3	40, 6±15, 5	42, 9±12, 7	36, 9±16, 5	49, 5±14, 1	47, 1±12, 5	
Averange	72, 6±25, 8	78, 8±25, 2	84, 4±14, 4	66, 8±22, 3	80, 7±20, 6	86, 6±11, 4	85, 7±11, 7	85, 2±12, 0	85, 5±13, 1	87, 3±10, 3	85, 2±11, 4
	66, 2±25, 9	71, 3±26, 0	78, 7±18, 0	65, 7±23, 1	57, 8±39, 2	80, 1±17, 3	79, 0±18, 4	81, 3±15, 1	80, 1±17, 9	71, 9±23, 9	79, 9±15, 9
	60, 1±25, 2	68, 2±25, 8	69, 0±26, 8	57, 0±28, 9	40, 4±14, 8	70, 3±26, 2	69, 3±25, 9	69, 8±25, 9	70, 6±26, 0	64, 3±29, 9	68, 7±25, 1
Wins	0/12	0/12	0/12	0/12	2/12	2/12	2/12	1/12	3/12	4/12	2/12
	1/12	0/12	1/12	0/12	0/12	4/12	2/12	3/12	2/12	1/12	3/12
	1/12	2/12	1/12	0/12	0/12	2/12	1/12	3/12	3/12	1/12	2/12

Tabla 12.1.: Resultados sobre las bases de datos del repositorio UCI (clasificación binaria). Las medidas Acc, GM y FM son consideradas. La media \pm la desviación estándar son mostradas como resultado de la validación cruzada anidada de 10-fold.

Datasets	<i>SV M</i>	<i>WSVM</i>	<i>SVM_{SMOTE}</i>	<i>WLTSVM</i>	<i>TBSVM</i>	<i>EWTSVM</i>	<i>EWTSVM*</i>	<i>EWTSVM**</i>	<i>EWTSVM[†]</i>	<i>EWTSVM^{†*}</i>	<i>EWTSVM^{†**}</i>
BankNote	7, 7±0, 3	7, 4±0, 7	59, 6±3, 2	102, 9±5, 2	217, 7±5, 5	215, 3±27, 9	127, 1±18, 7	135, 3±2, 5	170, 1±2, 4	49, 8±1, 2	92, 0±1, 1
Wisconsin	2, 1±0, 1	1, 8±0, 0	17, 6±0, 3	13, 0±0, 6	25, 0±0, 5	29, 9±0, 7	13, 6±0, 6	15, 0±0, 2	31, 2±0, 7	4, 3±0, 3	7, 7±0, 2
Ionosphere	1, 4±0, 1	1, 1±0, 0	11, 3±0, 3	5, 7±0, 3	10, 9±0, 1	14, 5±0, 3	5, 3±0, 0	8, 1±0, 2	14, 9±0, 1	1, 7±0, 1	3, 1±0, 1
Prima-Indians	12, 2±0, 5	12, 8±0, 6	128, 8±5, 2	22, 9±0, 9	52, 7±1, 0	57, 8±6, 2	17, 5±0, 2	25, 3±0, 2	50, 1±0, 7	7, 2±0, 3	13, 2±0, 4
Biodeg	9, 0±0, 3	8, 4±0, 5	89, 2±1, 5	47, 8±2, 2	119, 3±1, 8	104, 0±1, 0	39, 4±5, 0	58, 6±6, 9	104, 6±1, 1	13, 0±0, 3	24, 5±0, 5
Iris	0, 5±0, 1	0, 4±0, 0	3, 7±0, 0	1, 7±0, 0	6, 6±0, 1	5, 7±0, 0	1, 8±0, 1	2, 5±0, 1	6, 1±0, 1	0, 4±0, 0	1, 4±0, 1
HaberMan	6, 1±0, 6	3, 9±0, 3	49, 7±2, 6	4, 6±0, 6	10, 6±0, 4	17, 9±1, 4	3, 1±0, 2	5, 0±0, 4	12, 5±0, 1	0, 9±0, 2	3, 1±0, 1
Transfusion	14, 2±3, 8	13, 7±1, 7	79, 8±17, 0	20, 4±0, 7	53, 2±0, 6	52, 9±0, 7	11, 0±0, 8	33, 2±0, 8	53, 0±0, 8	2, 5±0, 1	35, 5±1, 0
Vehicle	3, 9±0, 1	3, 6±0, 1	43, 6±0, 6	27, 4±1, 0	72, 5±2, 3	68, 8±1, 0	14, 6±0, 4	48, 7±1, 8	69, 7±1, 0	3, 5±0, 2	50, 7±1, 5
cmc	22, 5±6, 3	26, 3±7, 5	377, 8±74, 9	92, 9±5, 2	317, 0±3, 6	256, 9±16, 4	54, 9±2, 2	203, 2±13, 7	249, 2±3, 3	11, 8±0, 5	189, 4±4, 3
Housing	2, 1±0, 1	2, 0±0, 1	42, 3±1, 8	8, 2±0, 5	27, 6±0, 6	30, 7±0, 9	3, 7±0, 1	20, 4±1, 2	31, 4±0, 8	0, 4±0, 0	25, 3±0, 8
Balance	8, 6±0, 8	10, 4±1, 8	149, 8±5, 9	15, 2±0, 7	38, 3±0, 4	42, 4±0, 9	5, 5±0, 1	31, 8±2, 0	42, 6±0, 8	0, 5±0, 0	40, 5±1, 5
Averange	7, 5±6, 4	7, 6±7, 5	87, 8±101, 7	30, 2±34, 0	79, 3±95, 7	74, 7±80, 5	24, 8±36, 0	48, 9±60, 4	69, 6±72, 7	8, 0±13, 9	40, 5±53, 6
Win	0/12	5/12	0/12	0/12	0/12	0/12	0/12	0/12	0/12	7/12	0/12

Tabla 12.2.: Tiempo de entrenamiento para las bases de datos del repositorio UCI (clasificación binaria). La media \pm desviación estándar son representan el tiempo de entrenamiento en segundos por cada fold.

ahí, podemos expresar que el $ETWSVM^{\dagger*}$ es la mejor opción con respecto a la relación al desempeño de la clasificación y el tiempo de entrenamiento.

También, se realizó una prueba estadística entre los clasificadores, realizamos la prueba de Kruskal-Wallis para comparar las medidas de desempeño entre los métodos. Si la hipótesis nula para la igualdad de medianas es rechazada, realizamos múltiples pruebas de comparación usando el algoritmo Tukey-Kramer para estudiar la diferencia entre los clasificadores [70, 71]. Todos los niveles de significancia con medidos en 5 %. Como se puede ver en la figura 12.6 SVM, WSVM, TBSVM, y WLTSVM exhiben comportamientos inestables a lo largo de las medidas Acc, GM y FM. Por otra parte, los algoritmos basados en ETWSVM y SVM_{SMOTE} son estadísticamente iguales. Aunque, el $ETWSVM^{\dagger*}$ muestra diferencias estadísticas en GM y FM en comparación a los demás algoritmos basados en ETWSVM y el SVM_{SMOTE} debido a la relación entre el rendimiento de la clasificación y el tiempo de entrenamiento, aún así se mantiene como una alternativa adecuada para lidiar con los desafíos del desbalance y traslape de los datos sin necesidad de conocimiento a priori en relación con el ajuste del algoritmo.

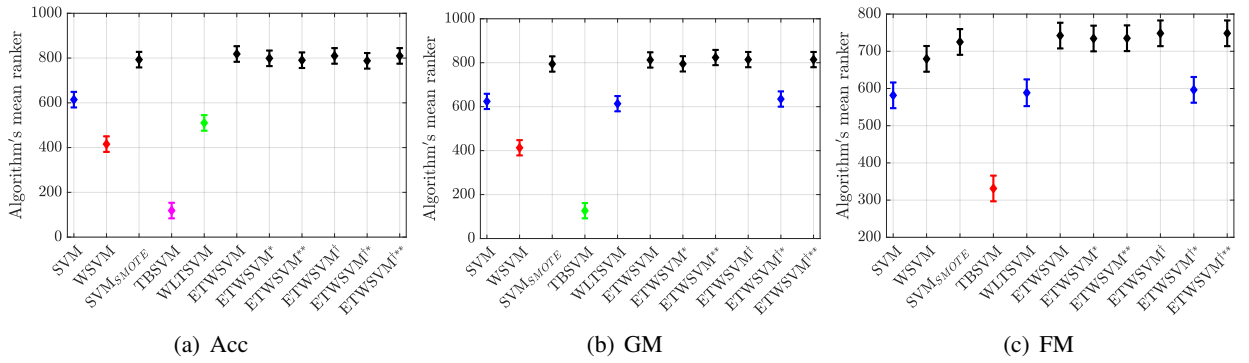


Figura 12.6.: Prueba de diferencia estadística entre los clasificadores basada en Tukey-Kramer. Se muestra el diagrama de caja del ranking de media del algoritmo, y el color indica que algoritmos son estadísticamente iguales.

12.3. Resultados bases de datos reales: clasificación multiclase

En la tabla 12.3 muestra el desempeño de la clasificación multiclase. Por un lado, SVM-OvR y SVM-OvO presentan resultados similares ($\overline{Acc} \sim 88\%$ y $FM_{\mu} \sim 78\%$). Por el otro lado, el TBSVM-OvO supera al algoritmo TBSVM-OvR; de hecho, el TBSVM-OvR presenta el más bajo desempeño entre los clasificadores estudiados. Esto último puede explicarse porque la estrategia OvR tiende a incrementar los problemas relacionados al desbalance de datos. Respecto a las variantes del ETWSVM, podemos ver como nuestras propuestas contrarresta, a través de un adecuado mapeo a un RKHS, los problemas de desbalance de datos independientemente de la estrategia multiclase (OvO o OvR). En general, el ETWSVM supera los métodos SVM y TBSVM, demostrando su capacidad de generalización de nuevo. En contraste con el problema de clasificación binario, el $ETWSVM^{\dagger*}$ no presenta el más bajo costo computacional por la necesidad de la estimación del kernel basado en CKA en varias ocasiones debido a los esquemas OvR y OvO.

Por último, realizamos las mismas pruebas estadísticas del escenario biclase. Como es evidente en la figura 12.7, los clasificadores basados en ETWSVM y el TBSVM-OvO presentan resultados consistentes en \overline{Acc} y FM_{μ} , mientras que, SVM-OvO, SVM-OvR, y TBSVM-OvR son estadísticamente diferentes de acuerdo a la prueba realizada.

	SVM OvR	SVM OvO	TBSVM OvR	TBSVM OvO	ETWSVM OvR	ETWSVM OvO	ETWSVM* OvR	ETWSVM* OvO	ETWSVM† OvR	ETWSVM† OvO	ETWSVM†* OvR	ETWSVM†* OvO
Datasets	Acc FM _μ	Acc FM _μ	Acc FM _μ	Acc FM _μ	Acc FM _μ	Acc FM _μ	Acc FM _μ	Acc FM _μ	Acc FM _μ	Acc FM _μ	Acc FM _μ	Acc FM _μ
Balance	94, 4±2, 3 91, 7±3, 4	94, 4±2, 3 91, 7±3, 4	38, 2±0, 9 7, 4±1, 3	99, 0±1, 2 98, 6±1, 8	98, 6±1, 2 97, 9±1, 9	98, 4±1, 3 97, 6±1, 9	98, 8±1, 2 98, 2±1, 8	98, 9±1, 0 98, 4±1, 5	98, 3±1, 5 97, 4±2, 3	98, 6±0, 9 97, 9±1, 3	97, 2±2, 4 95, 8±3, 6	99, 1±1, 3 98, 7±2, 0
Contraceptive	61, 8±0, 1 42, 7±0, 2	61, 8±0, 1 42, 7±0, 2	56, 7±0, 9 35, 1±1, 3	70, 0±2, 2 54, 9±3, 3	70, 3±2, 7 55, 4±4, 1	70, 2±2, 0 55, 3±3, 1	70, 0±2, 1 55, 1±3, 2	68, 7±3, 8 53, 0±5, 7	70, 1±3, 1 55, 2±4, 6	70, 1±2, 7 55, 2±4, 1	68, 8±2, 1 53, 2±3, 1	68, 3±3, 2 52, 5±4, 8
Dermatology	98, 6±0, 9 95, 9±2, 6	97, 5±1, 7 92, 6±5, 1	69, 6±1, 1 8, 8±3, 4	99, 1±0, 9 97, 3±2, 6	98, 6±0, 9 95, 9±2, 6	99, 1±0, 6 97, 3±1, 8	98, 8±1, 0 96, 4±2, 9	98, 9±0, 7 96, 7±2, 2	98, 7±0, 8 96, 2±2, 4	99, 3±0, 7 97, 8±2, 1	98, 2±0, 9 94, 5±2, 6	98, 7±0, 9 96, 2±2, 7
Ecoli	95, 3±1, 8 79, 0±7, 9	95, 2±1, 8 78, 3±7, 9	80, 3±2, 4 11, 5±10, 6	96, 0±0, 9 81, 9±3, 8	96, 8±1, 2 85, 8±5, 4	96, 7±1, 0 85, 1±4, 4	96, 5±1, 0 84, 3±4, 3	96, 5±1, 5 84, 0±6, 7	97, 0±1, 4 86, 7±6, 1	96, 9±1, 2 85, 8±5, 4	96, 5±0, 9 84, 3±4, 1	96, 6±1, 4 84, 6±6, 3
Glass	84, 9±2, 6 54, 6±7, 6	84, 4±2, 3 53, 2±6, 9	77, 3±1, 2 31, 9±3, 7	89, 2±3, 9 67, 5±11, 7	89, 1±3, 1 67, 3±9, 4	90, 0±4, 4 70, 1±13, 1	88, 2±3, 7 64, 5±11, 1	90, 1±2, 4 70, 2±7, 1	88, 8±3, 2 66, 4±9, 5	89, 7±3, 3 69, 1±9, 9	89, 0±4, 5 67, 0±13, 4	90, 6±3, 0 71, 9±9, 0
Hayes-roth	70, 2±5, 3 55, 3±7, 9	69, 7±5, 1 54, 6±7, 6	63, 2±5, 7 44, 7±8, 5	88, 9±7, 8 83, 3±11, 7	87, 9±6, 3 81, 9±9, 5	87, 8±6, 9 81, 7±10, 3	86, 5±7, 9 79, 7±11, 8	86, 9±4, 5 80, 4±6, 8	89, 4±7, 4 84, 1±11, 1	88, 0±4, 6 81, 9±6, 9	84, 4±7, 5 76, 7±11, 2	87, 8±5, 6 81, 7±8, 3
New-thyroid	97, 5±2, 5 96, 3±3, 7	97, 8±2, 1 96, 7±3, 1	57, 2±6, 9 35, 8±10, 4	96, 9±2, 5 95, 4±3, 8	97, 2±2, 8 95, 8±4, 2	96, 9±3, 0 95, 3±4, 4	90, 1±4, 6 85, 1±6, 8	97, 8±2, 6 96, 7±3, 8	97, 2±2, 2 95, 8±3, 4	97, 2±2, 3 95, 8±3, 4	98, 1±2, 2 97, 2±3, 3	97, 8±3, 0 96, 7±4, 5
Iris	81, 3±9, 8 72, 0±14, 7	81, 3±7, 5 72, 0±11, 2	55, 6±0, 0 33, 3±0, 0	96, 4±4, 1 94, 7±6, 1	97, 3±2, 3 96, 0±3, 4	97, 8±3, 1 96, 7±4, 7	93, 8±6, 4 90, 7±9, 5	94, 7±2, 8 92, 0±4, 2	96, 9±2, 2 95, 3±3, 2	97, 8±3, 1 96, 7±4, 7	93, 3±6, 0 90, 0±9, 0	94, 2±5, 2 91, 3±7, 7
Thyroid	95, 7±0, 9 93, 5±1, 3	95, 7±1, 1 93, 6±1, 6	43, 6±12, 2 15, 4±18, 3	97, 7±0, 9 96, 5±1, 3	97, 0±1, 5 95, 6±2, 3	96, 8±1, 3 95, 2±2, 0	95, 3±0, 7 92, 9±1, 0	97, 4±1, 5 96, 1±2, 3	96, 9±1, 7 95, 3±2, 5	97, 1±0, 9 95, 7±1, 4	97, 0±1, 6 95, 6±2, 4	97, 5±1, 4 96, 3±2, 1
Wine	97, 0±2, 9 95, 5±4, 4	96, 3±3, 0 94, 5±4, 5	51, 3±1, 3 27, 0±1, 9	98, 9±1, 8 98, 3±2, 7	98, 5±2, 6 97, 8±3, 9	98, 9±1, 8 98, 3±2, 7	98, 9±1, 8 98, 4±2, 6	97, 4±3, 9 96, 1±5, 9	99, 6±3, 1, 2 99, 4±1, 8	99, 2±2, 5 98, 8±3, 7	98, 5±1, 9 97, 8±2, 9	97, 4±2, 5 96, 1±3, 8
Penbased	97, 6±0, 5 87, 8±2, 7	97, 6±0, 6 87, 9±3, 2	81, 9±0, 1 9, 6±0, 6	99, 6±0, 3 97, 9±1, 5	99, 6±0, 2 97, 9±1, 0	99, 7±0, 2 98, 5±1, 2	99, 4±0, 2 96, 7±1, 1	99, 6±0, 3 98, 1±1, 3	99, 6±0, 3 97, 8±1, 5	99, 7±0, 3 98, 5±1, 5	99, 6±0, 2 97, 9±1, 1	99, 7±0, 2 98, 4±1, 0
Average	88, 6±12, 6 78, 6±19, 6	88, 4±12, 6 78, 0±19, 5	61, 4±14, 6 23, 7±13, 4	93, 8±8, 8 87, 8±14, 7	93, 7±8, 7 87, 9±14, 4	93, 8±8, 7 88, 3±14, 2	92, 4±8, 7 85, 6±14, 3	93, 4±9, 1 87, 4±14, 6	93, 9±8, 7 88, 2±14, 5	94, 0±8, 8 88, 5±14, 5	92, 8±9, 2 86, 4±14, 9	93, 4±9, 1 87, 7±14, 4
Wins	0/11 0/11	0/11 0/11	0/11 0/11	1/11 1/11	1/11 1/11	2/11 2/11	0/11 0/11	0/11 0/11	3/11 3/11	2/11 2/11	1/11 1/11	2/11 2/11

Tabla 12.3.: Resultados sobre las bases de datos del repositorio Keel (clasificación multiclase). Las medidas \bar{Acc} y el FM_{μ} son considerados. La media \pm desviación estándar son mostradas para la validación cruzada anidada de 10-fold.

Datasets	SVM OvR	SVM OvO	TBSVM OvR	TBSVM OvO	ETWSVM OvR	ETWSVM OvO	ETWSVM* OvR	ETWSVM* OvO	ETWSVM† OvR	ETWSVM† OvO	ETWSVM†* OvR	ETWSVM†* OvO
Balance	2, 2±0, 1	2, 0±0, 1	65, 4±1, 9	56, 3±1, 4	49, 4±0, 5	46, 5±0, 7	19, 3±1, 2	9, 4±0, 2	0, 8±0, 0	0, 7±0, 0	3, 0±0, 1	3, 1±0, 1
Contraceptive	169, 3±44, 3	157, 4±34, 2	500, 2±4, 7	292, 3±4, 9	302, 9±3, 4	202, 6±1, 6	243, 4±38, 7	166, 7±1, 9	4, 5±0, 2	2, 8±0, 1	78, 3±1, 9	80, 2±3, 2
Dermatology	3, 3±0, 2	6, 7±0, 0	65, 7±0, 9	83, 0±0, 5	55, 2±1, 0	75, 1±0, 5	15, 0±0, 4	12, 8±0, 2	0, 8±0, 0	1, 0±0, 0	3, 2±0, 1	3, 5±0, 2
Ecoli	9, 5±2, 3	15, 8±0, 4	85, 7±1, 2	164, 6±12, 3	72, 2±2, 8	145, 9±1, 4	12, 0±1, 1	15, 0±0, 2	1, 0±0, 0	1, 9±0, 1	0, 5±0, 1	0, 6±0, 1
Glass	10, 7±2, 2	7, 7±0, 5	30, 1±0, 3	70, 5±0, 3	25, 7±0, 3	58, 7±0, 2	5, 9±0, 1	7, 3±0, 1	0, 4±0, 0	0, 8±0, 0	0, 9±0, 0	0, 9±0, 0
Hayes-roth	1, 5±0, 0	1, 4±0, 0	8, 3±0, 1	13, 9±0, 1	7, 3±0, 0	11, 9±0, 0	3, 0±0, 1	2, 2±0, 1	0, 1±0, 0	0, 2±0, 0	0, 6±0, 0	0, 6±0, 0
New-thyroid	1, 8±0, 3	1, 3±0, 0	13, 4±0, 1	18, 4±0, 4	11, 3±0, 1	15, 6±0, 1	4, 0±0, 1	2, 6±0, 1	0, 2±0, 00	0, 2±0, 0	0, 7±0, 0	0, 7±0, 0
Iris	1, 7±0, 2	24, 3±0, 6	8, 9±0, 1	445, 0±5, 3	7, 8±0, 0	12, 0±0, 0	3, 2±0, 1	2, 4±0, 1	0, 1±0, 00	0, 2±0, 0	0, 7±0, 0	0, 7±0, 0
Thyroid	33, 8±32, 3	24, 2±8, 9	112, 3±1, 3	13, 6±0, 0	81, 9±0, 7	93, 0±1, 1	28, 6±1, 2	10, 1±0, 1	1, 1±0, 0	1, 2±0, 0	0, 6±0, 0	0, 6±0, 0
Wine	1, 4±0, 1	12, 5±5, 1	10, 8±0, 1	123, 4±2, 1	9, 4±0, 0	13, 5±0, 0	3, 3±0, 2	2, 9±0, 1	0, 1±0, 0	0, 2±0, 0	0, 9±0, 0	0, 9±0, 0
Penbased	152, 7±32, 1	1, 3±0, 0	1132, 7±13, 1	15, 5±0, 1	760, 7±9, 8	360, 5±5, 7	180, 5±17, 2	103, 3±1, 8	10, 5±0, 4	4, 8±0, 2	40, 7±0, 8	40, 6±1, 0
Average	35, 3±63, 0	23, 2±45, 4	184, 9±344, 6	117, 9±137, 5	125, 8±226, 9	94, 1±107, 4	47, 1±83, 1	30, 4±53, 8	1, 8±3, 2	1, 3±1, 4	11, 8±25, 1	12, 0±25, 5
Win	0/11	1/11	0/11	0/11	0/11	0/11	0/11	0/11	8/11	2/11	0/11	0/11

Tabla 12.4.: Tiempo de entrenamiento para el repositorio Keel (clasificación multiclase). La media \pm la desviación estándar mostradas corresponden al tiempo de entrenamiento y esta dado en segundos por fold.

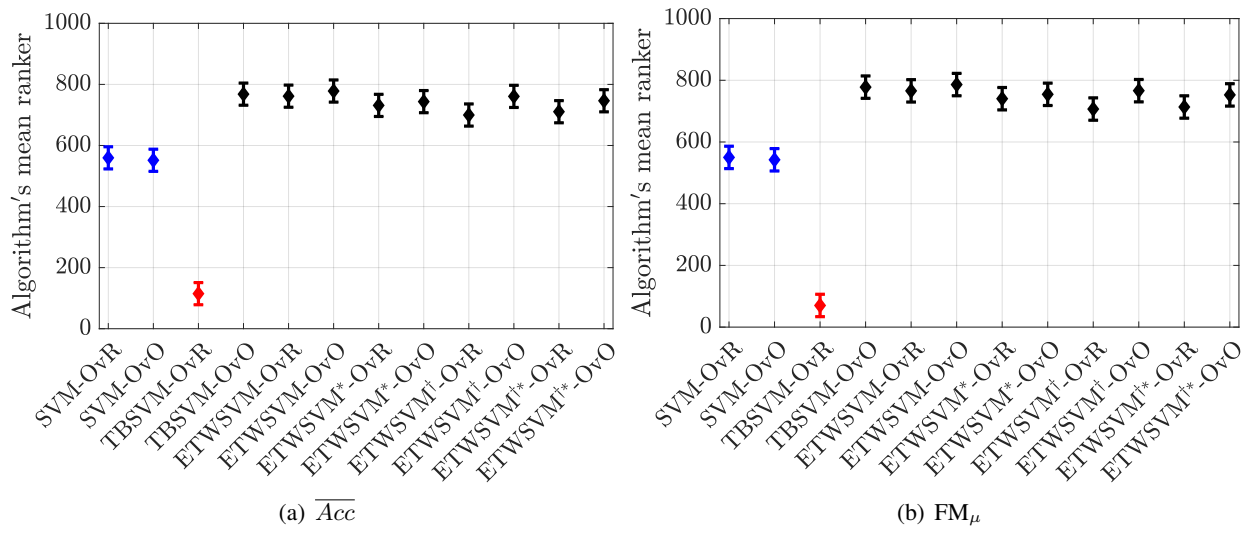


Figura 12.7.: Prueba de significancia estadística entre los clasificadores estudiados sobre el repositorio Keel, usando la prueba de Tukey-kramer. Se muestra el diagrama de caja del ranking de media del algoritmo, y el color indica que algoritmos son estadísticamente iguales.

13. Conclusiones

Presentamos un clasificador dedicado a la clasificación de datos desbalanceados nombrado *enhanced twin support vector machine*–(ETWSVM). Nuestra propuesta representa los datos de entrada en un espacio de características de alta dimensionalidad a través de un adecuado mapeo a un RKHS dentro de la optimización del ETWSVM. También, estimamos la función kernel usando una proyección lineal basada en CKA para emparejar las características de entrada con las etiquetas correspondientes. Entonces, el ETWSVM estimula la separabilidad de los datos y usa el conocimiento a priori para sintonizar la función kernel.

En este sentido, nuestro método es una nueva extensión no-lineal que incorpora una formulación dual en un RKHS de dimensión “infinita” para el clasificador TBSVM. Entonces, aplicamos el lemma de la matriz inversa para lograr una estimación basada en datos para los hiperplanos de separación. A su vez, el puntaje de clasificación es calculada como la distancia mínima entre los nuevos datos y los hiperplanos entrenados.

El ETWSVM básico es evaluado dada la configuración de los parámetros de regularización. Luego se estudio la influencia de los parámetros para diferentes escenarios de tasas de desbalance y traslape. Aparte de esto, el algoritmo SMOTE es acoplado con el ETWSVM por el bien de la comparación. Sucesivamente, extendemos el ETWSVM a problemas multiclase por medio de las estrategias OvR y OvO.

Los resultados experimentales se llevan a cabo sobre bases de datos sintéticas y reales. Específicamente, usamos los repositorios UCI y Keel para bases de datos reales para problemas biclase y multiclase, respectivamente. Los resultados logrados en términos del acierto, media geométrica, F-measure y tiempo requerido para el entrenamiento, muestran que nuestro ETWSVM y variantes superan a los métodos del estado del arte. Destacando, al ETWSVM^{†*} como una opción adecuada por su relación de desempeño de clasificación y tiempo de entrenamiento, sin necesidad de usar información a priori sofisticada para sintonizar los parámetros de regularización.

14. Trabajos futuros

Se planea extender el algoritmo ETWSVM para problemas de gran escala [72, 73]. También, probar varios tipos de funciones kernel con su estimación basada en CKA podría ser una línea de investigación interesante [74].

Con respecto al problema multiclase se piensa en la implementación de otras estrategias para extender el ETWSVM y sus variantes a este problema [39].

15. Publicaciones

Como resultado de esta investigación se realizaron las siguientes publicaciones:

Publicaciones en conferencias

- Jimenez C., Diaz D., Salazar D., Alvarez A.M., Orozco A., Henao O. (2018) Nerve Structure Segmentation from Ultrasound Images Using Random Under-Sampling and an SVM Classifier. In: Campilho A., Karray F., ter Haar Romeny B. (eds) Image Analysis and Recognition. ICIAR 2018. Lecture Notes in Computer Science, vol 10882. Springer, Cham
- Jimenez C., Alvarez A.M., Orozco A. (2019) A Data Representation Approach to Support Imbalanced Data Classification Based on TWSVM. In: Vera-Rodriguez R., Fierrez J., Morales A. (eds) Progress in Pattern Recognition, Image Analysis, Computer Vision, and Applications. CIARP 2018. Lecture Notes in Computer Science, vol 11401. Springer, Cham

Publicaciones en revista

- Jimenez C., Alvarez A.M., Orozco A. (2019), Imbalanced data classification using an enhanced twin support vector machine, In: Pattern Recognition. **(Sometido)**.

Bibliografía

- [1] S. Liu, Y. Wang, J. Zhang, C. Chen, and Y. Xiang. Addressing the class imbalance problem in twitter spam detection using ensemble learning. *Computers and Security*, 69:35–49, 2017. cited By 8. 1
- [2] M. Bach, A. Werner, J. Żywiec, and W. Pluskiewicz. The study of under- and over-sampling methods' utility in analysis of highly imbalanced data on osteoporosis. *Information Sciences*, 384:174–190, 2017. cited By 23. 1
- [3] S. Fotouhi, S. Asadi, and M.W. Kattan. A comprehensive data level analysis for cancer diagnosis on imbalanced data. *Journal of Biomedical Informatics*, 90, 2019. cited By 1. 1
- [4] L. Peng, H. Zhang, Y. Chen, and B. Yang. Imbalanced traffic identification using an imbalanced data gravitation-based classification model. *Computer Communications*, 102:177–189, 2017. cited By 8. 1
- [5] J. Sun, J. Lang, H. Fujita, and H. Li. Imbalanced enterprise credit evaluation with dte-sbd: Decision tree ensemble based on smote and bagging with differentiated sampling rates. *Information Sciences*, 425:76–91, 2018. cited By 24. 1
- [6] Guo Haixiang and et al. Learning from class-imbalanced data: Review of methods and applications. *Expert Syst. Appl.*, 73:220–239, 2017. 1, 20, 27
- [7] O. Loyola-González, J.F. Martínez-Trinidad, J.A. Carrasco-Ochoa, and M. García-Borroto. Study of the impact of resampling methods for contrast pattern based classifiers in imbalanced databases. *Neurocomputing*, 175:935–947, 2016. cited By 25. 1, 2, 42
- [8] P. Branco, L. Torgo, and R.P. Ribeiro. A survey of predictive modeling on imbalanced domains. *ACM Computing Surveys*, 49(2), 2016. cited By 82. 1, 2, 41
- [9] Lei Bao, Cao Juan, Jintao Li, and Yongdong Zhang. Boosted near-miss under-sampling on svm ensembles for concept detection in large-scale imbalanced datasets. *Neurocomputing*, 172:198–206, 2016. 1, 3
- [10] Maciej Zieba and Jakub M Tomczak. Boosted svm with active learning strategy for imbalanced data. *Soft Computing*, 19(12):3357–3368, 2015. 2, 3
- [11] Bernhard Schölkopf and Alexander J Smola. Learning with kernels: Support vector machines, regularization. *Optimization, and Beyond. MIT press*, 1(2), 2002. 2, 6, 8, 33
- [12] D. Tomar and S. Agarwal. Twin support vector machine: A review from 2007 to 2014. *Egyptian Informatics Journal*, 16(1):55–69, 2015. cited By 29. 2, 3
- [13] Z. Qi, Y. Tian, and Y. Shi. Structural twin support vector machine for classification. *Knowledge-Based Systems*, 43:74–81, 2013. cited By 97. 2

- [14] L. Cao and H. Shen. Combining re-sampling with twin support vector machine for imbalanced data classification. pages 325–329, 2017. cited By 3. 2
- [15] L. Liu, L. Wang, H. Ji, W. Zang, and D. Li. Between-class discriminant twin support vector machine for imbalanced data classification. volume 2017-January, pages 7117–7122, 2017. cited By 1. 2
- [16] Y. Xu, Q. Wang, X. Pang, and Y. Tian. Maximum margin of twin spheres machine with pinball loss for imbalanced data classification. *Applied Intelligence*, 48(1):23–34, 2018. cited By 4. 2
- [17] W. C. Lin, C. F. Tsai, Y. H. Hu, and J. S. Jhang. Clustering-based undersampling in class-imbalanced data. *Information Sciences*, 409-410:17–26, 2017. cited By 42. 2
- [18] N.V. Chawla, K.W. Bowyer, L.O. Hall, and W.P. Kegelmeyer. Smote: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*, 16:321–357, 2002. cited By 4583. 2, 42
- [19] S. Wojciechowski and S. Wilk. Difficulty factors and preprocessing in imbalanced data sets: An experimental study on artificial data. *Foundations of Computing and Decision Sciences*, 42(2):149–176, 2017. cited By 4. 2
- [20] Junru Lu, Chunkai Zhang, and Fengxing Shi. A classification method of imbalanced data base on pso algorithm. In *International Conference of Young Computer Scientists, Engineers and Educators*, pages 121–134. Springer, 2016. 2
- [21] Alejandro Moreo, Andrea Esuli, and Fabrizio Sebastiani. Distributional random oversampling for imbalanced text classification. In *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, pages 805–808. ACM, 2016. 2
- [22] Fulong Ren, Peng Cao, Wei Li, Dazhe Zhao, and Osmar Zaiane. Ensemble based adaptive over-sampling method for imbalanced data learning in computer aided detection of microaneurysm. *Computerized Medical Imaging and Graphics*, 55:54–67, 2017. 2
- [23] Jaesub Yun, Jihyun Ha, and Jong-Seok Lee. Automatic determination of neighborhood size in smote. In *Proceedings of the 10th International Conference on Ubiquitous Information Management and Communication*, page 100. ACM, 2016. 2
- [24] Chun-Wu Yeh, Der-Chiang Li, Liang-Sian Lin, and Tung-I Tsai. A learning approach with under-and over-sampling for imbalanced data sets. In *Advanced Applied Informatics (IIAI-AAI), 2016 5th IIAI International Congress on*, pages 725–729. IEEE, 2016. 2
- [25] X. Zhang, S. Ding, and T. Sun. Multi-class lstmsvm based on optimal directed acyclic graph and shuffled frog leaping algorithm. *International Journal of Machine Learning and Cybernetics*, 7(2):241–251, 2016. cited By 18. 2, 37, 39
- [26] W. Deng, L. Deng, J. Liu, and J. Qi. Sampling method based on improved c4.5 decision tree and its application in prediction of telecom customer churn. *International Journal of Information Technology and Management*, 18(1):93–109, 2019. cited By 0. 2
- [27] Yoav Freund, Robert E Schapire, et al. Experiments with a new boosting algorithm. In *Icml*, volume 96, pages 148–156, 1996. 3
- [28] X. Tao, Q. Li, W. Guo, C. Ren, C. Li, R. Liu, and J. Zou. Self-adaptive cost weights-based support vector machine cost-sensitive ensemble for imbalanced data classification. *Information Sciences*, 487:31–56, 2019. cited By 1. 3

- [29] M. Schlögl, R. Stütz, G. Laaha, and M. Melcher. A comparison of statistical learning methods for deriving determining factors of accident occurrence from an imbalanced high resolution dataset. *Accident Analysis and Prevention*, 127:134–149, 2019. cited By 0. 3
- [30] Yuan-Hai Shao, Wei-Jie Chen, Jing-Jing Zhang, Zhen Wang, and Nai-Yang Deng. An efficient weighted lagrangian twin support vector machine for imbalanced data classification. *Pattern Recognition*, 47(9):3158 – 3167, 2014. 3, 27, 40, 42
- [31] Fanyong Cheng, Jing Zhang, and Cuihong Wen. Cost-sensitive large margin distribution machine for classification of imbalanced data. *Pattern Recognition Letters*, 80:107–112, 2016. 3
- [32] Gerardo Casañola-Martin, Teresa Garrigues, Marival Bermejo, Isabel González-Álvarez, Nam Nguyen-Hai, Miguel Ángel Cabrera-Pérez, Huong Le-Thi-Thu, et al. Exploring different strategies for imbalanced adme data problem: case study on caco-2 permeability modeling. *Molecular diversity*, 20(1):93–109, 2016. 3
- [33] Sara Del Río, Victoria López, José Manuel Benítez, and Francisco Herrera. On the use of mapreduce for imbalanced big data using random forest. *Information Sciences*, 285:112–137, 2014. 3
- [34] V. López, S. Del Río, J.M. Benítez, and F. Herrera. Cost-sensitive linguistic fuzzy rule based classification systems under the mapreduce framework for imbalanced big data. *Fuzzy Sets and Systems*, 258:5–38, 2015. cited By 102. 3
- [35] Sang-Hoon Oh. Error back-propagation algorithm for classification of imbalanced data. *Neurocomputing*, 74(6):1058–1061, 2011. 3
- [36] D. Tomar and S. Agarwal. A comparison on multi-class classification methods based on least squares twin support vector machine. *Knowledge-Based Systems*, 81:131–147, 2015. cited By 65. 3, 39
- [37] A.J. Brockmeier, J.S. Choi, E.G. Kriminger, J.T. Francis, and J.C. Principe. Neural decoding with kernel-based metric learning. *Neural Computation*, 26(6):1080–1107, 2014. cited By 18. 3, 37
- [38] A.M. Alvarez-Meza, A. Orozco-Gutierrez, and G. Castellanos-Dominguez. Kernel-based relevance analysis with enhanced interpretability for detection of brain activity patterns. *Frontiers in Neuroscience*, 11(OCT), 2017. cited By 4. 3, 36, 40
- [39] S. Ding, X. Zhao, J. Zhang, X. Zhang, and Y. Xue. A review on multi-class twsvm. *Artificial Intelligence Review*, pages 1–27, 2017. cited By 1; Article in Press. 3, 33, 38, 42, 55
- [40] A. Roy, R.M.O. Cruz, R. Sabourin, and G.D.C. Cavalcanti. A study on combining dynamic selection and data preprocessing for imbalance learning. *Neurocomputing*, 286:179–192, 2018. cited By 8. 3
- [41] P. Xanthopoulos and T. Razzaghi. A weighted support vector machine method for control chart pattern recognition. *Computers and Industrial Engineering*, 70(1):134–149, 2014. cited By 42. 3
- [42] S. Piri, D. Delen, and T. Liu. A synthetic informative minority over-sampling (simo) algorithm leveraging support vector machine to enhance learning from imbalanced datasets. *Decision Support Systems*, 106:15–29, 2018. cited By 8. 3
- [43] Emanuel Parzen. Statistical inference on time series by hilbert space methods, i. Technical report, STANFORD UNIV CA APPLIED MATHEMATICS AND STATISTICS LABS, 1959. 6

- [44] Nachman Aronszajn. Theory of reproducing kernels. *Transactions of the American mathematical society*, 68(3):337–404, 1950. 6
- [45] Erwin Kreyszig. *Introductory functional analysis with applications*, volume 1. wiley New York, 1978. 6
- [46] George Kimeldorf and Grace Wahba. Some results on tchebycheffian spline functions. *Journal of mathematical analysis and applications*, 33(1):82–95, 1971. 8
- [47] William Karush. Minima of functions of several variables with inequalities as side constraints. Master’s thesis, Univ. of Chicago, 1939. 13
- [48] Harold W Kuhn and Albert W Tucker. Nonlinear programming. In *Proceedings of the 2nd Berkeley Symposium on Mathematical Statistics and Probability*, 1951. 13
- [49] Richard O Duda, Peter E Hart, and David G Stork. *Pattern classification*. John Wiley & Sons, INC, 2012. 14
- [50] Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012. 14
- [51] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, Sep 1995. 14
- [52] Yu-Xin Li, Yuan-Hai Shao, and Nai-Yang Deng. Improved prediction of palmitoylation sites using pwms and svm. *Protein and peptide letters*, 18(2):186–193, 2011. 15
- [53] Dino Isa, Lam H Lee, VP Kallimani, and Rajprasad Rajkumar. Text document preprocessing with the bayes formula for classification using the support vector machine. *IEEE Transactions on Knowledge and Data engineering*, 20(9):1264–1272, 2008. 15
- [54] Yingjie Tian, Yong Shi, and Xiaohui Liu. Recent advances on support vector machines research. *Technological and Economic Development of Economy*, 18(1):5–33, 2012. 15
- [55] Edgar Elias Osuna. *Support vector machines: Training and applications*. PhD thesis, Massachusetts Institute of Technology, 1998. 20
- [56] Olvi L Mangasarian and Edward W Wild. Multisurface proximal support vector machine classification via generalized eigenvalues. *IEEE transactions on pattern analysis and machine intelligence*, 28(1):69–74, 2006. 21
- [57] and R. Khemchandani and Chandra Suresh. Twin support vector machines for pattern classification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(5):905–910, May 2007. 22, 24, 25
- [58] Craig Saunders, Alexander Gammerman, and Volodya Vovk. Ridge regression learning algorithm in dual variables. 1998. 24
- [59] Y. Shao, C. Zhang, X. Wang, and N. Deng. Improvements on twin support vector machines. *IEEE Transactions on Neural Networks*, 22(6):962–968, June 2011. 26, 42
- [60] C. Cortes, M. Mohri, and A. Rostamizadeh. Algorithms for learning kernels based on centered alignment. *Journal of Machine Learning Research*, 13:795–828, 2012. cited By 151. 36

- [61] J. López, S. Maldonado, and M. Carrasco. A robust formulation for twin multiclass support vector machine. *Applied Intelligence*, 47(4):1031–1043, 2017. cited By 1. 38
- [62] Y.-H. Shao, W.-J. Chen, Z. Wang, C.-N. Li, and N.-Y. Deng. Weighted linear loss twin support vector machine for large-scale classification. *Knowledge-Based Systems*, 73(1):276–288, 2014. cited By 38. 39
- [63] X. Hua and S. Ding. Weighted least squares projection twin support vector machines with local information. *Neurocomputing*, 160:228–237, 2015. cited By 15. 39
- [64] M. Sokolova and G. Lapalme. A systematic analysis of performance measures for classification tasks. *Information Processing and Management*, 45(4):427–437, 2009. cited By 1109. 41
- [65] X.-Y. Liu, J. Wu, and Z.-H. Zhou. Exploratory undersampling for class-imbalance learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part B: Cybernetics*, 39(2):539–550, 2009. cited By 614. 42
- [66] X. Zhang, Y. Liang, J. Zhou, and Y. Zang. A novel bearing fault diagnosis model integrated permutation entropy, ensemble empirical mode decomposition and optimized svm. *Measurement: Journal of the International Measurement Confederation*, 69:164–179, 2015. cited By 163. 42
- [67] B.M. Abidine, L. Fergani, B. Fergani, and M. Oussalah. The joint use of sequence features combination and modified weighted svm for improving daily activity recognition. *Pattern Analysis and Applications*, 21(1):119–138, 2018. cited By 6. 42
- [68] J. Mathew, C.K. Pang, M. Luo, and W.H. Leong. Classification of imbalanced data by oversampling in kernel space of support vector machines. *IEEE Transactions on Neural Networks and Learning Systems*, 29(9):4065–4076, 2018. cited By 7. 42
- [69] José Alberto Hernández-Muriel, Andrés Marino Álvarez-Meza, Julián David Echeverry-Correa, Álvaro Ángel Orozco-Gutierrez, and Mauricio Alexander Álvarez-López. Feature relevance estimation for vibration-based condition monitoring of an internal combustion engine. *Tecno Lógicas*, 20(39):159–174, 2017. 48
- [70] J. Pizarro, E. Guerrero, and P.L. Galindo. Multiple comparison procedures applied to model selection. *Neurocomputing*, 48(1-4):155–173, 2002. cited By 40. 51
- [71] H.D. Vargas Cardona, A.A. Orozco, and M.A. Álvarez. Multi-patient learning increases accuracy for subthalamic nucleus identification in deep brain stimulation. *Conference proceedings : ... Annual International Conference of the IEEE Engineering in Medicine and Biology Society. IEEE Engineering in Medicine and Biology Society. Conference*, pages 4341–4344, 2012. cited By 1. 51
- [72] M. Goto, R. Ishida, and S. Uchida. A preselection-based fast support vector machine learning for large-scale pattern sets using compressed relative neighborhood graph. *Research Reports on Information Science and Electrical Engineering of Kyushu University*, 22(1):1–7, 2017. cited By 0. 55
- [73] S. Sharma and R. Rastogi. Stochastic conjugate gradient descent twin support vector machine for large scale pattern classification. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 11320 LNAI:590–602, 2018. cited By 0. 55

- [74] M. Zareapoor, P. Shamsolmoali, D. Kumar Jain, H. Wang, and J. Yang. Kernelized support vector machine with deep learning: An efficient approach for extreme multiclass dataset. *Pattern Recognition Letters*, 115:4–13, 2018. cited By 4. 55